

How do Categorical Duplicates Affect ML? A New Benchmark and Empirical Analyses

Vraj Shah[†]

IBM Research

vraj@ibm.com

Thomas Parashos

California State University, Northridge

thomasjparashos@gmail.com

Arun Kumar

University of California, San Diego

arunkk@eng.ucsd.edu

ABSTRACT

The tedious grunt work involved in data preparation (prep) before ML reduces ML user productivity. It is also a roadblock to industrial-scale cloud AutoML workflows that build ML models for millions of datasets. One important data prep step for ML is cleaning duplicates in the *Categorical* columns, e.g., deduplicating *CA* with *California* in a *State* column. However, how such *Categorical* duplicates impact ML is ill-understood as there exist almost no in-depth scientific studies to assess their significance. In this work, we take the first step towards empirically characterizing the impact of *Categorical* duplicates on ML classification with a three-pronged approach. We first study how *Categorical* duplicates exhibit themselves by creating a labeled dataset of 1262 *Categorical* columns. We then curate a downstream benchmark suite of 16 real-world datasets to make observations on the effect of *Categorical* duplicates on five popular classifiers and five encoding mechanisms. We finally use simulation studies to validate our observations. We find that Logistic Regression and *Similarity* encoding are more robust to *Categorical* duplicates than two *One-hot* encoded high-capacity classifiers. We provide actionable takeaways that can potentially help AutoML developers to build better platforms and ML practitioners to reduce grunt work. While some of the presented insights have remained folklore for practitioners, our work presents the first systematic scientific study to analyze the impact of *Categorical* duplicates on ML and put this on an empirically rigorous footing. Our work presents novel data artifacts and benchmarks, as well as novel empirical analyses to spur more research on this topic.

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/pvn25/CategDupsRepo>.

1 INTRODUCTION

Automated machine learning (AutoML) is beginning to increase access to ML for both small-medium enterprises and non-ML domain experts. This has led to the emergence of several platforms such as Google Cloud AutoML [7], Microsoft’s AutomatedML [10], and H2O Driverless AI [8] with the promise to automate the end-to-end ML workflow without any human-in-the-loop. Since ML prediction accuracy is the most critical in AutoML environments, many works have studied the automation and impact of algorithm selection, hyperparameter search, and optimization heuristics on ML [48, 53]. Also, recently there is a growing interest for studying how data prep specifically affects downstream ML [62, 69, 70].

Data prep for ML remains particularly challenging on structured data. It involves manual grunt work that is both tedious and

Table 1: A simplified dataset for ML Churn Prediction.

Name	Gender	State	Title	Contract	Zipcode	Density	MostCommonCrime_Zipcode	Churn
John	Male	California	sr. Scientist	monthly	93449	727	BURG	‘Y’
Jerry	Mail	CA	snr scientist	Month-to-month	91042	563	burglary	‘N’

time-consuming. Even AutoML users are often asked to manually perform many data prep steps before using their platforms [6]. Surveys of AutoML users have repeatedly identified such challenges in conducting data prep [42, 89]. One issue that they often encounter is duplicates in the columns that are *Categorical*, which assumes mutually exclusive values from a known finite set. This can require significant manual effort to fix duplicates even if a single *Categorical* column contains them in a data file.

Consider a dataset to be used for a common ML classification task in Table 1. Duplicates, categories referring to the same real-world object, occur in many *Categorical* columns such as *Gender*, *State*, *Title*, and *Contract*. Note that *Name* is not *Categorical* since it offers no discriminative power and cannot be generalized for ML. The presence of duplicates within a *Categorical* column can potentially dilute signal strength that one can extract for ML. Thus, an ML practitioner would often deduplicate categories before ML. We further discuss the conundrum for an ML practitioner in Section 3. Even, AutoML platforms often suggest users to manually inspect *Categoricals* and consolidate duplicates whenever they arise, as part of their guidelines for obtaining an accurate model [5]. This can involve non-trivial amount of deduplication effort at a *Categorical* column-level as duplicates can arise as misspellings, abbreviations, and synonyms, even within the same column. Note that this problem is related but complementary to entity deduplication issues studied in the data cleaning literature, as we explain in Section 2.1.

In this paper, we ask: *How do Categorical duplicates impact commonly used ML classifiers? Is category deduplication effort even worthwhile for ML? Is it always needed regardless of the employed Categorical encoding scheme?* We take a step towards answering these questions by developing an in-depth scientific understanding of the importance of category deduplication for ML classification (henceforth referred to as “ML”). Our objectives are two-fold. (1) Perform an extensive empirical study to measure the impact of *Categorical* duplicates on ML and distill the findings into actionable insights for handling them. This can help ML practitioners decide when and how to prioritise their cleaning effort. Moreover, this can enable AutoML platform builders design better ML workflows. (2) Present critical artifacts that can help advance the science of building AutoML platforms by providing researchers an apparatus to tackle open questions in this direction.

[†] Work done at IBM Research and University of California, San Diego.

Approach Overview. We identify that the impact on ML accuracy in presence of *Categorical* duplicates can be characterized with several explanatory variables (EVs) such as their duplication properties, training data properties, *Categorical* encoding, and ML model. Considering this, we make three-part contributions to cover our goals. (1) We produce labeled dataset to study how real-world *Categorical* duplicates arise. (2) We create a downstream benchmark suite to phenomenise the impact on ML on real-world data containing *Categorical* duplicates with multiple EVs. (3) Significance of each variable is hard to discern when all EVs act together. We use simulation study to disentangle the impact with each EV and explain the phenomenon discretely.

Empirical Evaluation. An empirical comparison of our downstream benchmark reveals that category deduplication can often improve the ML accuracy significantly, e.g., the median lifts in % accuracies due to category deduplication on *One-hot* encoded Logistic Regression (LR), Random Forest (RF), and artificial neural network (ANN) are 0.5, 1.5, and 2 (over 16 datasets) resp. Thus, LR gets impacted much less with *Categorical* duplicates than the high-capacity models. Overall, we make *eight* such observations on the significance of EVs with downstream benchmark. We validate them with simulation study and provide explanations into how ML models with different biases behave with *Categorical* duplicates.

Takeaways for Practitioners. We distill our empirical analysis into a handful of actionable takeaways for ML practitioners and AutoML developers. For instance, LR is more robust to the adverse impact of *Categorical* duplicates than high-capacity models as it overfits less. Also, *Similarity* encoding [39], Transformer-based embedding [64], and tabular representation learning-based method TABBIE [55] are more robust than other encodings to tolerate *Categorical* duplicates, thereby diminishing the utility of category deduplication task. We also expose a critical shortcoming of *One-hot* and *String* encoding [75], when *Categorical* duplicates arising in the deployment (or inference) but not during training can affect ML performance significantly.

Some of these insights may be considered folklore by practitioners, but this work is the first in-depth systematic scientific study to assess the impact of *Categorical* duplicates on ML. We explain the impact from the bias-variance tradeoff perspective to put empirical results on a rigorous footing. Our analyses can benefit practitioners to systematically understand the various EVs that matter for accuracy. Also, this can be useful to develop better practices and design ML workflows that are robust to *Categorical* duplicates. Moreover, our work opens up new research directions at the intersection of ML theory, data management, and ML system design.

Overall, our work is novel in terms of new labeled dataset, benchmark, and novel empirical analyses. We make four contributions:

1. **A new benchmark dataset.** To the best of our knowledge, this is the first work to curate a large labeled dataset specifically for *Categorical* duplicates where the entities are annotated. We present several insights that characterizes how *Categorical* duplicates exhibit themselves.
2. **Empirical benchmarking to understand the significance of category deduplication on ML.** Our curated downstream

benchmark containing “in-the-wild” datasets enables us to point out cases where the task may or may not benefit ML.

3. **Characterization of explanatory variables with simulation study.** Our study can disentangle and explain the impact of EVs on how *Categorical* duplicates affect ML.
4. **Utility of our work.** We present the first in-depth scientific empirical study to systematically characterize when and why category deduplication can help/not help ML. We present several practical insights for practitioners. We identify open questions for further research where our labeled data can be a key enabler to address them. Also, we open source our benchmark to enable more community-driven contributions [4].

2 RELATED WORK

2.1 Entity Matching (EM) and String Matching (SM) Approaches

EM, the task of identifying whether records from two tables refer to the same real-world entity has received much attention with rule-based [72, 78, 79], learning-based [58, 63, 67, 85, 90], semi-supervised [57], unsupervised [45, 88], active-learning [66] approaches, and even with Large Language Models (LLMs) [68, 87]. They operate at a tuple-level since they have access to the entire feature vectors of the two tables. Note that tuple-level duplicates do not necessarily imply duplication in *Categorical* strings, and also vice versa. Thus, the problem of EM is orthogonal to category deduplication. Admittedly, it is possible to view category deduplication as an extension of row-level deduplication but doing so is non-trivial. *Regardless, our focus is to study only the impact of category deduplication on ML and not how to perform deduplication or compare deduplication methods.* Thus, prior work on EM is complementary to ours in terms of utility for AutoML platforms.

SM, finding strings from two sets that refer to the same real-world entity has been explored with an active learning solution Smurf [38] and an unsupervised learning approach [61]. However, such SM approaches are orthogonal to our focus on studying how *Categorical* duplicates affect ML. We leave automating category deduplication to future work, including potentially extending existing row-level deduplication and SM works.

2.2 Data Cleaning and Data Prep for ML

CleanML [62] analyses the impact of many data cleaning steps on ML. Our work is along the same direction, but they do not specifically explore *Categorical* deduplication and its causal variables that matter for accuracy. Although they do study string-level inconsistencies within a column with four real datasets, they are not all *Categorical*. Also, they focus on deriving a broad perspective and a coarse-grained study of many cleaning steps such as this. In contrast, we dive deep into *Categorical* deduplication. We study the different explanatory variables that matter for accuracy to offer empirical rigor and understand the importance of task scientifically.

We performed an objective benchmarking of a specific ML data prep step, namely the feature type inference task [77]. We build upon our open-sourced datasets but we study a completely different problem. There exist numerous data prep tools such as rule-based [54], exploratory data analysis-based [74], program synthesis [52, 56], and visual interfaces [13] to reduce manual grunt work

effort and allow users to productively prepare their data for ML. Our work’s insights can complement all these tools to reduce human time and effort and make their analysis more interpretable. Some works have studied human-in-the-loop cleaning to improve ML accuracy and reduce user effort [59, 60]. However, they do not support a cleaning operation with *Categorical* duplicates. Our labeled data can spur more follow-up works in this general direction of automating and improving data prep for ML. Error detection [41], ML for data cleaning methods [76, 83], and even techniques that perform value standardization [36] are orthogonal to our focus since we do not propose new techniques for *Categorical* deduplication.

2.3 AutoML Platforms

Several AutoML tools such as AutoML Tables [7], TransmogrifAI [12], and AutoGluon [46] performing automated model selection do automate many data prep tasks. However, they do not explicitly handle *Categorical* duplicates. Instead, the users are asked to explicitly clean and remove inconsistencies in *Categorical* columns before using their platform [6]. Our labeled data can lead to contributions from community to automate deduplication with a supervised learning-based approach, including potentially extending AutoML with deduplication processor in the optimization process [47, 69, 70]. Moreover, we believe that our empirical analyses and takeaways provide valuable insights to improve AutoML platforms.

3 OUR APPROACH

Consider again an ML practitioner predicting Customer Churn with Table 1 data. She sources the data from multiple tables such as *Customers* and *Zipcode*, which contains details about the area where the customers live. She expects likely duplicates in many *Categorical* columns such as *Gender*, *State*, *Title*, and *Contract* since they are collected using “Free Text” customer surveys. She wants to build and periodically retrain ML pipelines such that they are most robust to likely duplicates. She prefers an ML pipeline which is not necessarily the most accurate but the one that is most reliable. She wants to minimize any adverse impact from *Categorical* duplicates.

To build such an ML pipeline, she wants to choose from different *Categorical* encoding schemes and ML models that are popular on tabular data [80, 82]. Moreover, she has an intuition that many *Categoricals* such as *MostCommonCrime_Zipcode* are not relevant for the target and cleaning them may not be in the best interest of her time. She would like to prioritize her efforts towards cleaning *Categoricals* that are more likely to impact ML. In addition, she has the resources to collect even more training data on customers, but doesn’t know if more training data would necessarily translate to a more robust pipeline. Overall, she is fraught with several questions: *How do Categorical duplicates impact the behavior of popular ML classifiers and encodings? Would the effort towards cleaning duplicates be less worthwhile for a non-relevant column as opposed to a column that is relevant for the task? Would collecting more training data help in mitigating the impact of Categorical duplicates?*

Towards answering these questions, we take a step towards empirically assessing the significance of category deduplication on ML. We first identify the important variables that matter for ML practitioners and study how they affect ML. We hand label a large dataset of real-world *Categorical* columns with duplicates to

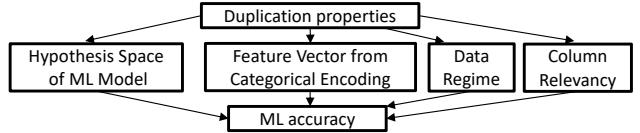


Figure 1: Summarization of EVs impacting ML in the context of *Categorical* column that has duplicates.

understand how they occur. We then make empirical observations of the impact of deduplication with different explanatory variables (EVs) in the real-world. We finally use synthetic study to validate observed phenomenon and intricately study how each EV impact ML. We first summarize the EV we study and then explain our three-part contributions towards building an in-depth understanding of the importance of category deduplication for ML.

We focus this study in the context of a *Categorical* column that has duplicates. As the domain size of column shrinks with deduplication, it can influence the following EVs impacting ML (as Figure 1 shows): (1) Feature vector from *Categorical* encoding. (2) Hypothesis space denoting a set of all prediction functions from feature space to label space that the ML model can represent. (3) Data regime in terms of the number of training examples per unique category in the column. (4) Column relevancy as a measure of the importance of column for the downstream task. Admittedly, there can exist other complex EVs such as skew in class labels with different distributions and conditional duplication properties given the class label. However, in this work, we focus on the above EVs because of their importance for ML practitioners. We leave studying the other EVs to future work.

1. Our Hand-Labeled Data. We create the first large labeled data where true entities within a *Categorical* column are annotated with duplicate categories. This helps us understand the observed properties of *Categorical* duplicates and how they manifest themselves in real-world columns. Our data includes 1262 string *Categorical* columns from 231 raw CSV files. The labeling process took us about 150 man-hours across 6 months. The utility of our labeled data is two-fold. (1) Configure duplication parameter ranges and skew distributions in simulation study. (2) Presents a crucial artifact for researchers to automate the task of *Categorical* deduplication itself and even to objectively evaluate the accuracy of in-house automated mechanisms by AutoML platform developers. In fact, one such labeled data for ML feature type inference task lead to objective benchmarking of existing AutoML tools and even lead to more accurate supervised ML approaches to automate the task [77]. We dive into this part in Section 5.

2. Downstream Benchmark Suite. We use 16 real-world datasets to empirically study the impact of *Categorical* duplicates. We choose these datasets such that they capture different kinds of duplication and also represents the different regimes in the EV’s spectrum (explained in Section 6.2). We choose five popular ML classifiers and five *Categorical* encoding schemes to showcase how they impact the behavior of duplicates on ML. We leave in-depth discussion of this component in Section 6.

3. Synthetic Study. We perform a Monte Carlo-style simulation study to achieve two objectives. (1) Confirm the validity of observations we make with downstream benchmark. (2) Disentangle and

Table 2: Notations used in this paper with a simplified example to illustrate our notions with *State* column categories.

Symbol	Meaning			
C	Set of category values in the column A_l			
E	Set of unique real-world entities referred by categories from C			
ED	Subset of real-world entities that have at least 1 duplicate; $ED \subseteq E$			
$\text{occ}(Z)$	Sum of occurrences of all categories present in set Z; $Z \subseteq C$			
D	Set of non-empty sets of duplicate values for each entity in ED; $ D = ED $			
Category set C_i ($1 \leq i \leq C $)		Occurrence of Category ($\text{occ}(\{C_i\})$)	Entity set E_j ($1 \leq j \leq E $)	
New York	C_1	60	New York	E_1
NY	C_2	30		
new york	C_3	10		
California	C_4	70	California	E_2
Ca	C_5	30		
Wisconsin	C_6	100	Wisconsin	E_3

characterize the effect of duplicates with multiple EVs individually to make the impact interpretable. We embed a true distribution and vary the EV one at a time while fixing the rest to study their impact on ML along with how they trend. Although we use our labeled data to inform duplication parameter values, our simulation study is not entirely dependent on it. One can very well fix arbitrary duplication parameter values, although that doesn't change the trends and conclusions we derive. Section 7 covers this in depth.

4 PRELIMINARIES

4.1 Scope

4.1.1 Focus on studying downstream ML. We focus on the ML classification setting over tabular data. We call the ML model to be trained over the data as the “downstream model.” We study how *Categorical* duplicates impact the performance of popular downstream models used ML practitioners [81, 82]. Note that our focus is not to study how *Categorical* duplicates interact with different optimization schemes for model selection or AutoML procedures [71, 84]. We leave this study to future work. Our goal is not to study the upstream deduplication process itself, which is handled manually in the paper. Also, our focus is not to extend existing entity matching approaches such as with LLMs [68, 87] or devise new methods for category deduplication. We leave studying automated upstream deduplication mechanisms to future work.

4.1.2 Focus on duplicates in *Categorical* columns. We focus on studying duplicates in the context of string *nominal Categorical* features, which do not have a notion of ordering among its values. Note that a *Categorical* feature contains mutually exclusive values from a known finite domain set. In contrast, *Text* type features can take arbitrary string values. Thus, generic open domain addresses or even non-generalizable person names are not used as *Categorical*. We study duplicates arising in *Categorical* feature column, which is not the actual target for the prediction task. We leave studying duplicates arising in the target column to future work.

4.2 Definitions

We present terms and notations needed to study the effect of *Categorical* duplicates in the context of implications for ML accuracy. We first draw upon notations from a mix of both database theory [65] and ML literature [51] for known concepts. A relational table is defined by schema $R(A_1, A_2, \dots, A_n, Y)$ with a relation (instance) r . We

use \mathcal{A} to denote a set of columns $\{A_1, A_2, \dots, A_n\}$ and Y is the target column for prediction. Note that, formally, a column is referred to as an attribute [65]. Let $A_l (l \in [1, n])$ be a *Categorical* column with a domain $\text{dom}(A_l) \subseteq \mathcal{L}$, where \mathcal{L} is the set of strings with finite length. A relation r is defined over \mathcal{A} as a set of mappings with $\{t^p : \mathcal{A} \rightarrow \bigcup_{l=1}^n \text{dom}(A_l), p = 1 \dots |r|\}$, where for each tuple $t^p \in r$, $t^p(A_l) \in \text{dom}(A_l)$, $|r|$ is the number of examples in the the table.

Note that *Categorical* strings are not directly consumable by most ML models. Thus, an encoding scheme is required to transform \mathcal{A} to a feature vector to train an ML model. We explain this further in Section 6.1. We now reuse and adapt terminologies from existing database [40, 65] and ML literature [51] together for terms that we need for the rest of the paper. Table 2 lists the notations and explains the terms used with an example. For simplicity of exposition, we focus on one *Categorical* column with duplicates, $A_l \in \mathcal{A}$.

DEFINITION (CATEGORY). A Category set $C^l = \{C_1^l, C_2^l, \dots, C_{|C^l|}^l\}$ contains all unique domain values occurring in the column A_l . Note that C^l is also referred to as the active domain of A_l relative to relation r [65], i.e., $C^l = \text{adom}(A_l, r) = \{c \in \text{dom}(A_l) \mid \exists t^p \in r, t^p(A_l) = c\}$. We drop the superscript (C^l) and simplify the active domain operation with C only to make it succinct for follow up set algebra. Each distinct value in the column is defined as “category.” For Table 2 example, $C = \{\text{New York}, \text{NY}, \text{new york}, \text{California}, \text{Ca}, \text{Wisconsin}\}$.

DEFINITION (ENTITY). An Entity set $E \subseteq C$ represents a subset of *Categories* that conceptually refer to different real-world objects. A category from set C can be uniquely mapped to an entity from set E . Let the mapping function be denoted by $M : C \rightarrow E$. In Table 2, there are three unique real-world state objects, i.e., $E = \{\text{New York}, \text{California}, \text{Wisconsin}\}$. Note that entities are defined at a conceptual level; thus, referring to New York as new York or NY is identical. But for ease of exposition, we assume the category that most frequently represents an entity (ties broken lexicographically) in the column to be the true entity. There exist multiple categories representing the same entity, i.e., $M(C_1) = M(C_2) = M(C_3) = E_1 = \{\text{New York}\}$.

DEFINITION (OCCURRENCE). We define Occurrence (or percentage Occurrence) of category C_i as percentage of times C_i represents E_j in the column. For instance, whenever real-world New York entity occurs, 30% and 10% of the times NY and new york represents them respectively. New York is referred to as the entity since it occurs more than NY and new york. We define the Occurrence function as $\text{occ} : Z \rightarrow [0, 100]$. The input Z is a subset $Z \subseteq C$ such that all categories of the subset map to a unique entity $E_j (j \in [1, |E|])$, i.e., $E_j = M(Z_1) = M(Z_2) = \dots = M(Z_{|Z|})$. The output is the sum of occurrence values for all categories present in the input set which is a real number in $[0, 100]$. $\text{occ}(Z) = \text{occ}(Z_1) + \dots + \text{occ}(Z_{|Z|})$, e.g., $\text{occ}(\{C_1\}) = 60$, $\text{occ}(\{C_2, C_3\}) = 40$, and $\text{occ}(\{C_1, C_4\}) = \text{Undefined}$.

DEFINITION (DUPLICATE). There exist a duplicate for E_j whenever $E_j = M(Z_1) = M(Z_2) = \dots = M(Z_{|Z|})$, $|Z| > 1$. Whenever E_j occurs, the % times it is represented by Z_1, Z_2 , and Z_n are $\text{occ}(Z_1), \text{occ}(Z_2)$, and $\text{occ}(Z_n)$ resp. Without loss of generality, we assume that $\text{occ}(Z_1) \geq \text{occ}(Z_2) \geq \dots \geq \text{occ}(Z_{|Z|})$. Since Z_1 most frequently represents the entity (ties broken lexicographically), the other categories Z_2, \dots, Z_n are referred to as duplicates of the entity E_j . We define $ED \subseteq E$ as the subset of the entities that contain at least one duplicate, i.e.,

Table 3: Duplication types w/ examples from our labeled data

Duplication Types	Column name	Category Examples
1 Capitalization	Country	"United States", "united States"
2 Misspellings	Gender	"Male", "Mail", "Make", "msle"
3 Abbreviation	State	"California", "CA"
	preparer_title	"Senior Counsel", "Sr. Counsel"
4 Difference of Special Characters	City	"New York", " New York, "
	Colour	"Black/Blue", "Black-Blue"
5 Different Ordering	Colour	"GoldWhite", "WhiteGold"
6 Synonyms	Gender	"Female", "Woman"
	Venue	"Festival Theatre", "Festival Theater"
7 Presence of Extra Information	City	"Houston", "Houston TX", "Houston TX 77055"
	Colour	"triColor", "tricolored"
8 Different grammar	Venue	"Auditorium", "TheAuditorium"

$\exists Z \subseteq C$ s.t. $|Z| > 1$ and $M(Z_1) = \dots = M(Z_{|Z|}) = ED_j$ ($j \in [1, |ED|]$). We define a duplicate set D_k ($k \in [1, |ED|]$) for every entity in ED such that $D_k = \{Z_2, Z_3, \dots, Z_{|Z|}\}$ represents a set of duplicate values, e.g., $ED_1 = \{California\}$, $D_1 = \{Ca\}$ and $ED_2 = \{New York\}$, $D_2 = \{new york, NY\}$.

DEFINITION (CATEGORY DEDUPLICATION). This is the task of mapping categories from C to an entity from E with mapping function M . The new column after the assignment is called the *deduplicated* column. Set C and E of the *deduplicated* column are identical.

DEFINITION (COLUMN RELEVANCY). Let $Acc(\mathcal{A})$ be the % classification accuracy obtained by the ML model with a set of columns \mathcal{A} to be used as features in the input. *Relevancy* of a column $A_l \in \mathcal{A}$ is defined as $Acc(\mathcal{A}) - Acc(\mathcal{A} - \{A_l\})$. This quantifies the absolute predictive power of column A_l for the downstream task.

5 OUR HAND-LABELED DATASET

We create a labeled dataset of *Categorical* columns where *Entities* in each column is annotated with their duplicates whenever present. This enables us to understand how real-world duplicates manifest themselves and what do the sets E, ED, D and their occurrences look like. We now discuss how this dataset is created, the types of real-world duplicates present, and our dataset analysis with stats and important insights into the behavior of duplicates.

5.1 Data Sources

We constructed a large real-world data of 9921 columns sourced from Kaggle and UCI ML repository with diverse application domains such as retail, healthcare, and finance [77]. Columns were manually annotated with a standardized 9-class vocabulary of ML feature types. The classes include feature types such as *Numeric*, *Categorical*, *Sentence*, and *Not-Generalizable* (e.g., primary keys). Using this, we obtain just the string *Categorical* columns. In addition, we collect more data files with such columns using open-source data portals [14, 16, 18, 21, 23, 25, 28]. We use 16 data files for empirical benchmarking on real downstream tasks in Section 6. Overall, we find 231 raw CSV files with at least one string *Categorical* column. We find a total of 1262 such columns.

Current Limitation. We sourced many *Categorical* columns by leveraging our previous dataset [77]. The raw files were collected from sources such as Kaggle and UCI ML repo where the data file may have been subjected to some pre-processing. However, this is the

best we can do from academic research standpoint given legal constraints: acquire large public datasets using public APIs, annotate them, and make them available to the community. It is hard to acquire truly raw data files from several enterprises and make them public due to legal constraints. Also, we do not make any general claims about the manifestation of duplicates across the universe of the datasets. This would require doing a comprehensive analysis of datasets from all sources including that from enterprises and other organizations. However, this does not diminish the utility of our empirical analyses as both the downstream benchmark suite and our synthetic study are independently useful.

5.2 Labeling Process

Among the *Categorical* columns we collected, we do not know which columns contain duplicates beforehand. This necessitates us to manually scan through all the 1262 *Categorical* columns and look for duplicates in them. We follow the below process at a column-level to reduce the cognitive load of labeling. For every *Categorical* column, we enumerate its category set with the count of times each category appears in it. Before scanning the category set, we sort the categories by their appearance count in descending order and their values in lexicographic order. As we scan the category set, we annotate duplicates with their corresponding entities in the column. We construct E, ED , and D sets, along with their occurrences for all the columns. *The entire labeling process took us roughly 150 man-hours across 6 months and 3 people.*

5.3 Types of Duplicates and Data Statistics

We find that there exist *eight* types of duplication. We present these types with examples in Table 3. The differences shown are relative to the representation of the true entity. We now clarify some of the types. *Type 4* denotes the difference of any non-alphanumeric special characters including comma, period, and white spaces. *Type 5* denotes different ordering within multi-valued categories. *Type 8* categories have either a common stem/lemma, presence of stop-words, or a common singular representation. Note that a duplicate can have duplication of multiple types and an entity can have numerous duplicates, each belonging to multiple types, e.g., given $ED_1 = \{New York\}$ and $D_1 = \{new-york, NY\}$, "new-york." has both *Type 1* and *4* duplication, and the entity *New York* has duplicates with duplication of *Type 1, 3, and 4*.

We annotated 67060 entities across all 1262 string *Categorical* columns. We find that almost 5% of those entities have the presence of at least one duplicate with a total of 5584 duplicates. Overall, 66 columns from 47 raw CSV files have the presence of at least one duplicate. There are three parameters that quantify the amount of duplication within a column. (1) Fraction of entities that have at least one duplicate ($|ED|/|E|$). (2) Duplicate set size for all entities of the column (set D). (3) Duplicate occurrences $occ(\{D_{ki}\})$, $k \in [1, |ED|], i \in [1, |D|]$. We present the complete stats on the different parameters that characterizes duplication over our labeled data and on duplication types in Appendix.

6 DOWNSTREAM BENCHMARK

We now empirically study the impact of category duplicates on the downstream ML tasks. Note that our focus is not to compare

Table 4: Statistics of the column containing *Categorical* duplicates in our 16 downstream datasets. $|r|$, $|\mathcal{A}|$, and $|Y|$ are the total number of examples, columns, and target classes in the data resp. Duplication types numbering correspond to Table 3. $|rC|$ denotes the number of training examples per category of set C . We use colors green, blue, red with hand-picked thresholds to visually present and better interpret the cases where the amount of duplication is low ($1-|E|/|C| < 0.25$), moderate ($1-|E|/|C| > 0.25 \& < 0.5$), and high ($1-|E|/|C| > 0.5$) resp. We use following thresholds with the same colors to better interpret the data regime: low ($|rC| < 5$), moderate ($|rC| > 5 \& < 25$), and high ($|rC| > 25$). Note that data regime moves up with category deduplication as category set size has shrunk. We present more fine-grained stats in the Appendix.

Datasets	$ r $	$ \mathcal{A} $	$ Y $	Duplication Types								$ C $	$1 - \frac{ E }{ C } \%$	Row $ rC $	Truth $ rC $
				1	2	3	4	5	6	7	8				
Midwest Survey	2778	29	9	X	X	X	X	X	X	X	X	1008	64	2.5	6.5
Mental Health	1260	27	5	X	X	X			X	X		49	69	23.2	81.2
Relocated Vehicles	3263	20	4	X	X	X	X			X		1097	36	2.5	3.8
Health Sciences	238	101	4	X		X			X			56	61	3.6	8.3
Salaries	1655	18	8		X		X			X		647	29	1.8	2.2
TSM Habitat	2823	48	19	X			X		X	X		912	11	2.6	2.9
EU IT	1253	23	5	X	X	X	X	X	X			256	35	3.9	5.9
Halloween	292	55	6	X			X			X	X	163	51	1.5	3
Utility	4574	13	95	X			X			X		199	31	16.2	24.3
Mid or Feed	1006	78	5	X		X	X		X			37	62	20.6	59.7
Wifi	98	9	2	X			X			X		69	52	1.3	2.5
Etailing	439	44	5	X	X	X	X			X		71	68	5.3	14.3
San Francisco	148654	13	2	X		X						2159	10	46.3	50.9
Building Violations	22012	17	6		X	X				X		270	63	53.7	145
US Labor	210287	25	4	X	X	X	X	X	X		X	1169	47	31	60.4
Pet Registration	82545	14	2	X			X	X		X		789	44	58.2	105

and evaluate category deduplication methods. We curate a benchmark suite of 16 real-world datasets, each containing a column with duplicates. We use this to empirically evaluate and compare three *Categorical* encoding schemes both with and without the presence of duplicates. Finally, we make several important observations on the different EVs that impact the relationship of *Categorical* duplicates with downstream classifiers.

6.1 Models and Encodings

We choose five classifiers used commonly among the ML practitioners [81, 82]: LR, RF, XGBoost (XGB), SVM using a Radial Basis Function (RBF) kernel (henceforth referred to as SVM), and an ANN. RF, SVM, and XGB offers large VC dimensions and are high-capacity, while LR has a low-capacity. We use ANN architecture with 2 hidden units, each with 100 neurons. Although there is no magic number for ANN architecture size, the above network already offers a very large VC-dimension and a high-capacity [37]. We use the synthetic study (Section 7) with two extremes in the ANN’s bias spectrum to empirically assess the different capacities of ANN on ML.

We encode *Categorical* columns with five popular schemes: (1) *One-hot (OHE)*, (2) *String (StrE)* [75], (3) *Similarity (SimE)* [39], (4) a pre-trained Transformer-based embeddings (*TransE*) [64], and (5) table embedding method *TABBIE* [55]. *OHE* is the standard approach to encode nominal *Categoricals* as it follows their two

properties. (1) Each category is orthogonal to one another. (2) Pairwise distance between any two categories is identical. With a category set C^l (for A_l) closed during training, *OHE* sets feature vector $X_l^p = [1(t^p(A_l)=C_1^l), \dots, 1(t^p(A_l)=C_{|C^l|}^l)]$, where $1(\cdot)$ is the indicator function and $p=1..|r|$. RF with *OHE* performs binary splits on the data. RF can also handle raw “stringified” *Categorical* values by performing set-based splits on the data. We refer to this as *StrE*. Note that *StrE* is not applicable for models which are not tree-based, since they cannot handle raw string values. Both *OHE* and *StrE* assume that the *Categorical* domain is closed with ML inference, i.e., new categories in the test not seen during training are handled by mapping them to a special category, “*Others*.”

SimE takes into account the morphological variations between the categories. The feature vector for category set C^l is given as $X_l^p = [Sim(t^p(A_l), C_1^l), \dots, Sim(t^p(A_l), C_{|C^l|}^l)]$, where $Sim(\cdot)$ is a similarity metric defined as the dice-coefficient over n -gram (n ranges from 2 to 4) strings [35]. *TransE* uses a pre-trained BERT base transformer model to obtain embeddings as features [64]. The feature vector can be computed even for any new categories arising in test set which are unseen during training for both *TransE* and *SimE*. *TABBIE* [55] is a pre-trained tabular model that consists of three transformers for converting rows, columns, and cells to vector representations. It first obtains an uncontextualized embedding for cells using the BERT model. The cell embeddings are then passed through row/column transformers, which see the rows/columns of the table as inputs and produce contextualized output representation for rows/columns respectively. The row/column embeddings are obtained by prepending CLSROW/CLSCOL tokens and using their embedding as the embedding for the full row/column. For downstream ML classification, since we want a prediction class from the target column and for every row in the test set, we first obtain the row-level representations given by *TABBIE*. We then place a single-layered ANN classifier (ANN-1L) output layer over *TABBIE* and we finally fine-tune the representations over the training data. We do not add more neural net layers to exclusively observe the effects of duplicates with *TABBIE* on ML.

6.2 Datasets used for Analyses

We choose 16 datasets from Section 5.1 such that we not only represent the different duplication types but also span the spectrum of different variable combinations. Table 4 presents the statistics over our datasets. We use the quantity % reduction in domain size with deduplication ($1-|E|/|C|$) to summarize the magnitude of duplication. We use the data regime notion to denote the number of training examples per category value of the column with duplicates ($|rC|$). We ensure that our selected datasets sufficiently represent different ranges of values (high vs. low measured relatively) in both EV spectrum. For instance, a dataset that involves a high amount of duplication coupled with high- and low-data regimes such as *Building Violation* and *Midwest Survey* respectively. We will later see in Section 6.4 that the former dataset is robust to duplicates even with almost 51% of their column’s entity diluted with duplicates, while the latter is not. This enables us to make specific observations on the role of different EVs, which we validate and disentangle using our simulation study in the Section 7.

We obtain the following data files: Midwest Survey[26], EU IT[31], Mental Health[29], Wifi[33], Relocated Vehicles[20], Utility[22],

Table 5: Classification accuracy comparison of ML models with different *Categorical* encodings over 16 downstream datasets. Accuracy results are shown as delta lift or drop (shown as negative values) in % diagonal accuracy with *Truth* (column with *Categorical* duplicates has been deduplicated with truth) relative to *Raw* (column has *Categorical* duplicates intact). Green, blue, and red colors denote cases where the *Truth* accuracy relative to *Raw* is significantly higher, comparable, and significantly lower (error tolerance of 1%) respectively. *TRel* denotes the true Relevancy of the column that has been deduplicated; this indicates the importance of the column for downstream target. Due to space constraints, we defer the absolute accuracy results on *Raw* and *Truth* with each method to Appendix.

Dataset	Random Forest				XGBoost			ANN			SVM		Logistic Regression		TABBIE + ANN-1L
	<i>TRel</i>	<i>OHE</i>	<i>StrE</i>	<i>SimE</i>	<i>OHE</i>	<i>StrE</i>	<i>SimE</i>	<i>OHE</i>	<i>SimE</i>	<i>TransE</i>	<i>OHE</i>	<i>SimE</i>	<i>OHE</i>	<i>SimE</i>	
Midwest Survey	16.1	11.5	10	4.4	9.8	10.1	3.4	9.5	3.8	8.5	9.4	0.9	9.4	2.1	7.4
Mental Health	1.3	1.1	-0.1	-1.7	1.5	0.8	2	2	-0.4	-0.7	1.6	0.1	1.3	0.6	1.2
Relocated Vehicles	9.1	3	4.1	-0.1	5.9	7.5	0.3	3.6	0	1.6	4.7	-0.2	4	0.4	1.9
Health Sciences	0.4	2.2	0	-2.7	0.4	-0.4	0.9	4.9	1.8	0.4	1.9	1.2	0.9	1.8	0.1
Salaries	0.7	1.7	1.3	0.4	-0.3	0.7	0.2	0.5	5.4	3.8	0	0	0.2	-1.3	4.5
TSM Habitat	5.2	0.4	1.4	0.4	0.9	2.1	-0.3	-2.7	-2.7	0	0.2	0.2	0	0	0.7
EU IT	3.3	1.2	-0.6	4	2.4	-1.1	0.6	-2.4	5	1.5	2.5	0.8	0	0	2.9
Halloween	-0.4	1.5	1.5	-4.9	2.3	4.2	0	4.2	0.8	0	3	3.8	3.4	1.1	0.1
Utility	8.1	1.4	1.2	1.4	0	1.1	-0.2	2.3	2.5	-0.2	0	0.5	-0.2	0.3	1.1
Mid or Feed	1.5	2.5	-0.2	1.8	3.3	0.3	0	2	0.2	0.1	0.2	0.3	1.7	-1.2	-0.4
Wifi	4.2	5.3	4.2	3.2	5.3	0	3.2	2.1	3.2	-0.9	1.1	8.4	1.1	8.4	0.3
Etailing	-0.5	2	1.1	3	-0.9	2.3	-0.7	-3	0	-0.7	0.6	0.5	-0.5	1.8	0.2
San Francisco	24.4	0.1	-0.3	0	0	-0.1	-0.1	0.1	-0.1	0.2	0.1	0.2	-0.1	0	0.1
Building Violations	-0.1	-0.1	0.1	0	0	0	0	0	0	-0.6	0	0	0	0	0
US Labor	3.9	1.3	0.8	1	1	0.5	1.1	2	0.8	0.7	2.1	1.4	0.8	0.4	0.8
Pet Registration	1.8	0.2	0.2	-0.1	0.2	0	0.1	0.2	0.2	0.2	0.2	0.1	0.2	0.1	-0.1

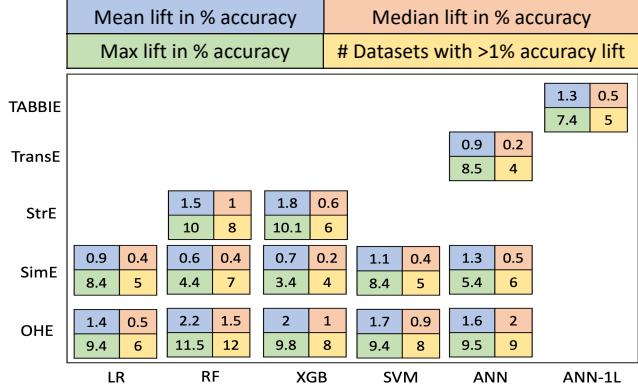


Figure 2: Summary statistics showing the impact on ML with *Raw* (all *Categorical* duplicates are intact) relative to *Truth*.

Health Sciences[17], Salaries[24], TSM Habitat[15], Etailing[34], Building Violations [19], Mid or Feed[32], Halloween[27], San Francisco [30], US Labor[1], and Pet Registration[2]. Each dataset has a column with *Categorical* duplicates which we manually deduplicated in Section 5.2. We do not claim that these 16 datasets are representative of the percentage one encounters in practice. Our goal with the downstream benchmark is not to make universal claims about the impact of *Categorical* duplicates on just the commonly encountered datasets. Instead, we select them plainly to showcase different EV settings and study the behavior of duplicates in those settings. The benchmark suite helps us point out the cases where deduplication matters. This coupled with synthetic study only serves as a guide to help ML practitioners and AutoML developers glean insights. We hope our work inspires more data benchmark standardization in this space with industry involvement.

6.3 Methodology

We partition each dataset into an 80:20 split of train and test sets. We perform 5-fold cross-validation and use a fourth of the examples in the train set for hyper-parameter search. We tune the regularization parameter for LR and SVM. We tune the number of trees and their maximum depth for RF and XGB with values for each ranging from 5 to 100. ANN is L2 regularized and tuned. We present the grids for hyper-parameter tuning in Appendix.

6.4 Results

6.4.1 Results comparing the ML impact with and without *Categorical* duplicates. Table 5 shows the delta lift/drop in diagonal accuracy of all ML models built with different encoding schemes when the downstream data is completely deduplicated with ground truth. As an example, *Raw Midwest Survey* contains a column with *Categorical*s duplicates. Cleaning its duplicates with their true entities (*Truth*) leads to an 11.5% lift in accuracy relative to the *Raw*, when building RF with *OHE*. Figure 2 shows statistics summarizing the impact with all ML models using different encodings on 16 downstream datasets. Finally, we present the generalization performance of classifiers with the overfitting gap (difference between train and validation accuracies) in Table 6. We summarize our results with important observations below.

O1. There exist several downstream cases where *Truth* improves the ML accuracy over *Raw* for any encoding scheme. For instance, the delta accuracy increase with *Truth* on RF with *OHE* is of median 1.5% and up to 11.5% compared to *Raw* (across 16 datasets). Moreover, the delta accuracy increase is of median 2% and up to 9.5% for ANN.

O2. Delta increases in accuracies with *Truth* are typically higher with all high-capacity models (RF, SVM, XGB, and ANN) than a low-capacity model LR. The median delta increases in accuracy with RF

Table 6: Comparisons of overfitting gap with three classifiers presenting representative choices from the bias-variance tradeoff spectrum using *OHE*. The drop in overfitting gap for *Truth* is shown relative to the *Raw*.

Dataset	RF		ANN		LR	
	Raw	Truth	Raw	Truth	Raw	Truth
Midwest Survey	50.7	-14.2	45.1	-10.4	24.4	-9.4
Mental Health	42.3	-7.2	26.7	-0.2	11.7	-3.5
Relocated Vehicles	27.3	-3.1	16.4	-3.6	17	-4.1
Health Sciences	35	-8.1	44.9	-4.9	9.3	-5.9
Salaries	34.6	-1	1.4	-0.5	1.9	+0.2
TSM Habitat	28	-0	0.1	+0.5	1.9	-0
EU IT	53.1	-6.6	1.4	+0.9	1.2	-0
Halloween	50.9	-5.8	58.1	-4.2	38.3	-3.5
Utility	41.2	-1.4	26.1	-3	0.7	-0.3
Mid or Feed	58.4	-1.1	66	-2	34.2	-12.8
Wifi	26.2	+1.3	47.4	-2.1	11.1	-2.1
Etailing	54.4	-1.6	59.7	+2.9	41.2	-7.7
San Francisco	-0.2	-0	1.1	-0.1	0.5	-0
Building Violations	1.8	-0.1	1.1	-0.2	0.2	+0.1
US Labor	20.8	-2.7	20.1	-3.1	1.6	-0.4
Pet Registration	1.1	+0.1	1.2	-0.3	1.1	-0.2

and ANN using *OHE* are 1.5 and 2, compared to 0.5 for LR. Thus, LR is more robust to duplicates than the high-capacity models.

O3. Categorical duplicates impact *OHE* the most, *StrE* (for RF and XGB) the second most, and encoding methods such as *SimE*, *TABBIE*, and *TransE* (for ANN) are the least affected (see Figure 2). Interestingly, the median lifts in accuracies due to deduplication with *SimE* are less than 0.5 for all models. Overall, *Truth* helps to improve the ML performance with *SimE* significantly in just $\sim 34\%$ of the all the 80 downstream cases (16 datasets and 5 models). This is because, *SimE* considers morphological variations between the category strings and maps a duplicate to a similar feature vector as the true entity. So, duplicates are often located close to their true entities in the feature space. Thus, any further lift in accuracy due to deduplication is marginal. Moreover, the embedding-based methods such as *TransE* and *TABBIE* exhibit marginal impact with duplicates, as the median lifts with *Truth* are 0.2 and 0.5 resp. *TABBIE* does not further help *TransE* in reducing the adverse impact of duplicates.

O4. The overfitting gap reduces with deduplication (from Table 6), thereby improving their generalization ability. Since high-capacity models such as RF and ANN are more prone to overfitting than LR, their accuracy lifts with *Truth* are more significant.

O5. If the magnitude of overfitting gap on *Raw* is insignificant ($< 1\%$), the amount of possible reduction in overfitting with *Truth* is also small. Thus, it's not worthwhile to deduplicate if the overfitting gap on *Raw* is already low to begin with. We observe this will all the datasets where the overfitting gap is close to 1%, e.g., *San Francisco* and *Building Violations*. We observe this trend across all classifiers.

O6. Category deduplication increases the column *Relevancy* for all models, i.e., the column becomes more predictive for the downstream tasks after category deduplication. Note that the magnitude of accuracy lift with *Truth* quantifies the increase in column *Relevancy* with *Truth*, as per definition in Section 4.2.

O7. The accuracy lifts with *Truth* on all the models are more significant when the column has high *Relevancy* unless there exist

a high-data regime with a large number of training examples per category. Thus, if a column has already high *Relevancy* on *Raw*, it may be worthwhile conservatively to deduplicate, e.g., *Relocated Vehicles* and *Midwest Survey*.

O8. High-data regime is robust to the impact of *Categorical* duplicates than low-data regime, regardless of the amount of duplication. Even a high amount of duplication has a negligible impact in the high-data regime, e.g., *Building Violations* has a massive 63% reduction in domain size due to deduplication, but there exist a large number of training examples per category. We do not see any lift in accuracy with category deduplication on any of the ML models.

Results with additional evaluation metrics. We rerun our downstream benchmark with metrics such as macro/micro average of precision, recall, and F1-score. *We find that none of the empirical conclusions made with diagonal accuracy change with these metrics.* We defer their results and discussion to Appendix.

6.4.2 Results with different duplication types. We now empirically benchmark the impact on ML with each duplication type discretely (from Table 3). Recall that we know the duplication type for every duplicate arising in our labeled data from Section 5. We utilize this to (1) pick a duplication type (say, Type k , $k \in [1, 8]$) (2) pick a *Raw* downstream dataset where Type k occurs, (3) deduplicate and consolidate duplicates of Type m ($\forall m \in [1, 8], m \neq k$) with their corresponding entities (categories that most frequently represent the true real-world concepts in the column; see definition in Section 4.2) in the *Raw* data. We then obtain a version of *Raw* (say, *Raw* _{k}) where only a single duplication type (Type k) is present. We repeat steps 2 and 3 for every *Raw* data in our downstream suite where Type k is present. We repeat all three steps $\forall k$ and obtain a suite of downstream datasets \mathcal{D}_k for every Type k , $k \in [1, 8]$.

We re-run the ML models together with encoding schemes from Section 6.1 on the full data suite \mathcal{D}_k , $k \in [1, 8]$. Figure 3 (A) shows the results summarizing how each duplication type impacts ML under different model/encoding settings. To better interpret our findings, Figure 3 (B) presents string similarity scores to understand how close the feature vector representations of duplicates are relative to their true entities. Dice coefficient score [35] counts the number of overlapping n -grams (n from 2 to 4) relative to the sum of total n -grams between the duplicate string and its corresponding true entity string. Euclidean measure calculates the distance in a pre-trained transformer-based embedding space [64] between the embeddings of duplicate relative to its true entity. We summarize our findings by the feature representation method used below.

OHE and *StrE*. Both methods model categories as mutually exclusive and they do not get affected by the semantics of *Categorical* strings. For instance, *OHE* masks the string with a binary variable, one for each unique category in the column. *StrE* perform set-based splits on the *Categorical* domain. Thus, we observe a similar trend as Section 6.4, across the duplication types: high-capacity models exhibit more adverse impact to duplicates than the low-capacity LR with *OHE* and *StrE* is more robust to duplicates than *OHE*.

SimE. We find this method to have good robustness to all duplicate types when used in conjunction with LR, except on *Type 7*. In contrast, high-capacity models manifest significant vulnerability when there exist duplicates such as *Type 3*, *Type 6*, and *Type 7* (Abbreviations, Synonyms, and Having Extra Information resp.).

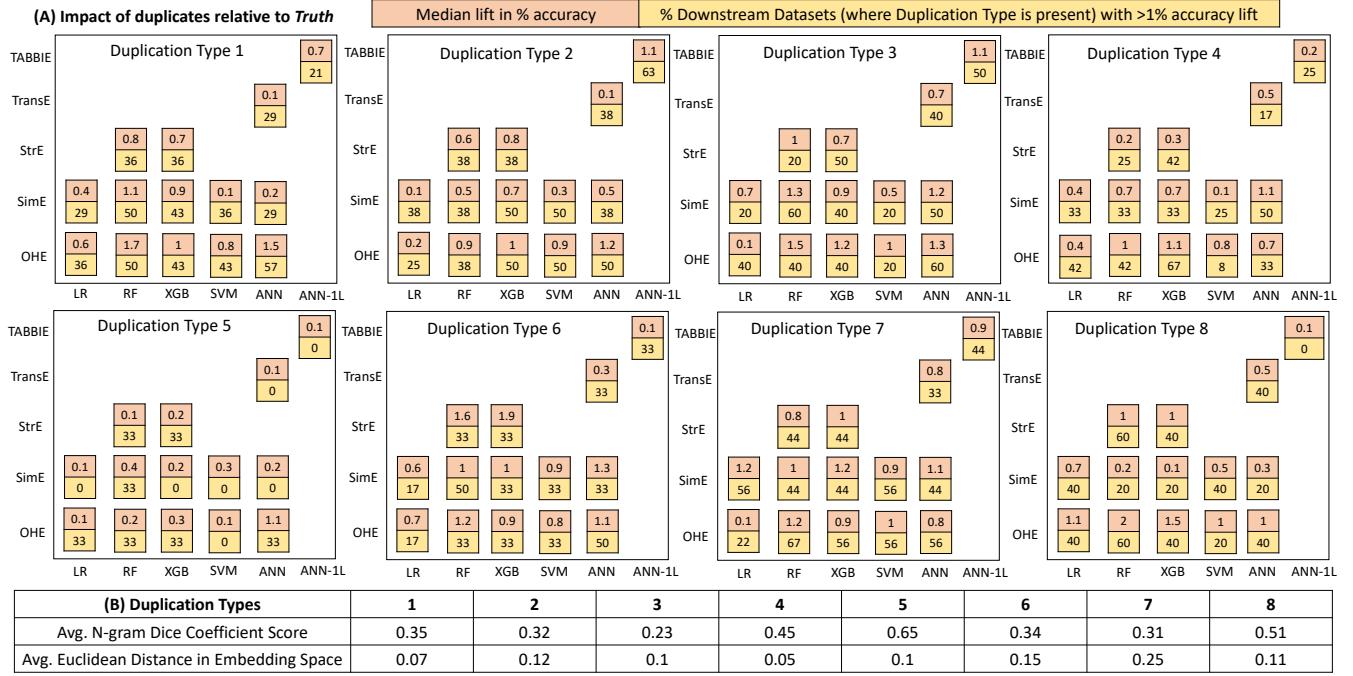


Figure 3: (A) Results showing the impact on ML exclusively with each duplication type from Table 3. Number of downstream datasets (Table 4 shows their names with types) corresponding to each duplication type are: Type 1) 14, Type 2) 8, Type 3) 10, Type 4) 12, Type 5) 3, Type 6) 6, Type 7) 9, Type 8) 5. This amounts to a total of 67 downstream cases across the 8 duplication types. **(B)** String similarity scores of duplicate w.r.t its true entity: Dice coefficient score [35] and Euclidean measure in the transformer embedding space [64]. The scores are computed for all duplicates present in our 16 downstream datasets and they are averaged across a given Type. A higher dice coefficient score and a lower euclidean distance imply that the feature representations of duplicates are closer to their true entities.

Interestingly, the dice-coefficient scores for such duplication types are also lower (relative to other Types), implying the feature vectors of duplicates to be far from their true entities. Moreover, duplicates due to Capitalization differences (*Type 1*) particularly cause tree-based models, RF and XGB to have a significant impact, as their median accuracy lifts with *Truth* and the number of datasets with significant performance differences are significantly high.

TransE and *TABBIIE*. We find both methods to be highly robust to almost all duplication types, except that *TABBIIE* exhibits significant performance impact with Misspellings and Abbreviations (*Type 2* and *Type 3* resp.). Although *Type 7* on *TransE* has the highest (compared to the other types) Euclidean distance (in the embedding space) from duplicates to their entities, the median lift with deduplication is not significantly high. The total number of downstream cases (out of 67) where we notice a $> 1\%$ performance impact of duplicates with *TransE* and *TABBIIE* are 20 and 21 resp. Thus, interestingly, learning representations over tables doesn't significantly help to subside the adverse impact of *Categoricals* over *TransE*.

7 IN-DEPTH SIMULATION STUDY

We now dive deeper into the impact of each EV on the downstream ML. This study helps us not only validate our empirical observations but also makes the significance of each EV impacting ML more interpretable. Moreover, it reveals the limitations of commonly used encoding schemes when unseen duplicates during training arise in the test.

7.1 Models and Encodings

The structural model parameters such as the number of tree estimators and maximum tree depth for RF and the specific ANN architecture can largely impact the bias-variance tradeoff. Thus, we fix them to disentangle their impact and better illustrate our findings by presenting two extremes of RF's and ANN's bias spectrum. We use high-bias models such as shallow decision tree with a restricted tree depth of 5 (denoted as ShallowDT), a low-capacity ANN comprising of two hidden units with 5 neurons each (denoted as LoCapANN), and also LR. In addition, we use low-bias high-capacity RF with the number of tree estimators and maximum tree depth being fixed to 50 (denoted as HiCapRF). These values represent the median best-fit parameters obtained by performing a grid search (with the grids being same as Section 6.3) over the synthetically generated data described in Section 7.2. We again use a high-capacity ANN comprising of two hidden units with 100 neurons each (HiCapANN). We also use LR, SVM, and XGB (using the same methodology as Section 6.3).

We focus this study in the context of *OHE* and *StrE*. *SimE* and *TransformE* require the categories to be semantically meaningful strings. An entity can have duplication of multiple types. Constructing a fine-grained simulator that generates semantically meaningful duplicates while preserving the same true entity is non-trivial and intricate from the language standpoint. We leave designing an apt simulation mechanism for *SimE* and *TransformE* to future.

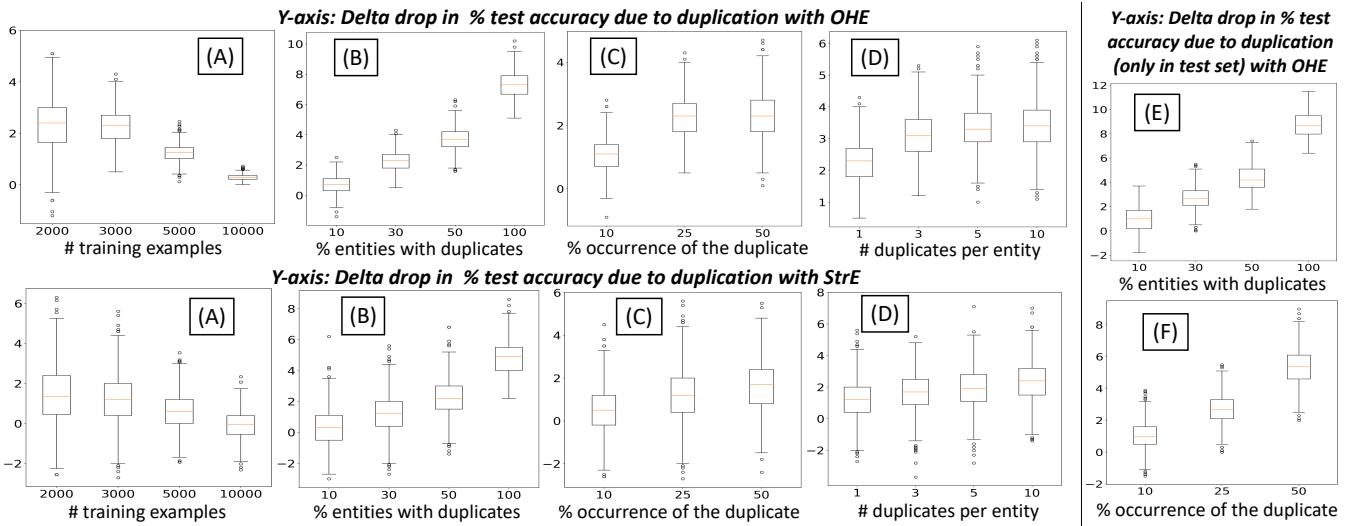


Figure 4: Simulation results for HiCapRF with OHE and StrE. (A-D) Duplicates are present in train, validation, and test set. (E-F) Only test set is diluted with duplicates. (A) Vary $|r|_t$ (# training examples) while fixing $(|ED|/|E|, occ(D_k), |D_k|)=(30, 25, 1)$ (B) Vary $|ED|/|E|$ while fixing $(|r|_t, occ(D_k), |D_k|)=(3000, 25, 1)$ (C) Vary $occ(D_k)$ while fixing $(|r|_t, |ED|/|E|, |D_k|)=(3000, 30, 1)$ (D) Vary $|D_k|$ while fixing $(|ED|/|E|, |r|_t, occ(D_k))=(30, 3000, 25)$, for all $k \in [1, |ED|]$. Parameter settings of (E) & (F) are same as (B) & (C) resp.

7.2 Setup and Data Synthesis

There is one relational table with Y being boolean. We include 3 *Categorical* columns in the table and set $|\mathcal{A}|=3$. We set entity set size of every column, $|E|=10$ (all columns have a domain size of 10).

Data generating process. We set up a “true” distribution $P(\mathcal{A}, Y)$ and sample examples in an IID manner. We study a complex joint distribution where all features obtained from \mathcal{A} determine Y . We sample $|r|$ number of total examples, where the examples for training, validation, and test are in 60:20:20 ratio. We then introduce synthetic duplicates in one of the columns of the table in different ways. We vary the EVs one at a time and study their impact on ML accuracy along with how they trend. We generate 100 different (clean) training datasets and 10 different dirty datasets for every clean one. We measure the average test accuracy and the average overfitting gap of all models obtained from these 1000 runs.

The exact sampling process is as follows. (1) Construct a conditional probability table (CPT) with entries for all possible values of \mathcal{A} from 1 to $|E|$. We then assign $P(Y=0|\mathcal{A})$ to either 0 or 1 with a random coin toss. (2) Construct $|r|$ tuples of \mathcal{A} by sampling values uniform randomly from $|E|$. (3) We assign Y values to tuples of \mathcal{A} by looking up into their respective CPT entry. (4) We perform the train, validation, and test split of this clean dataset and obtain the binary classification accuracy of the ML models on the test split.

Duplication process. We introduce duplicates in a column $A_l \in \mathcal{A}$ of the clean data as follows. (1) Fix fraction of entities to be diluted with duplicates, e.g., $|ED|/|E|=0.3$ (2) Form set ED (entities to be diluted with duplicates) by sampling uniformly randomly $|ED|$ categories from E , e.g., $ED=\{E_3, E_5, E_8\}$. (3) For every entity in ED , fix duplicate set size $|D_k|, k \in [1, |ED|]$, e.g., $|D_k|=1, k \in [1, 3]$. We assume that all entities have identical duplicate set sizes. (4) Given $|D_k|$, we form the set D by introducing duplicates, e.g., $D_1=\{E_3\text{-duplicate}_1\}, D_2=\{E_5\text{-duplicate}_1\}, D_3=\{E_8\text{-duplicate}_1\}$. (5)

Fix $occ(D_k), k \in [1, |ED|]$. For every duplicate value d in D , set occurrence $occ(d)=occ(D_k)/|D_k|$, i.e., assume that all the duplicates representing an entity are equally likely to occur. (6) We perform the same train, validation, and test split of the resulting dataset as obtained in step 4 of the data generating process. We finally obtain the test accuracy of the ML models on the dirty dataset. We use our labeled data to configure apt duplication parameter values such that we can showcase an average and worst-case scenario.

7.3 Results

We vary all EVs one at a time while fixing the rest. We confirm the trends and observations made with *italics*.

7.3.1 Varying the data regime. Figure 14 (A) presents the delta drop in %accuracy with duplication relative to the ground truth on HiCapRF as the number of training examples ($|r|_t$) are varied with both *OHE* and *StrE*. We find that with the rise in $|r|_t$, the delta drop in accuracy decreases. With just 3 training examples per CPT entry ($|r|_t = 3k$ and total entries in CPT=1k), duplicates cause a drop of median 2.3% and up to 4.3% accuracy with *OHE*. With 10 training examples, the median and max drops in accuracies due to duplicates with *OHE* are 0.3% and 0.7% respectively. This confirms our observation on the downstream benchmark suite: *A higher data regime is more robust to duplication than a lower data regime. The same trend holds with StrE encoding and also all the other classifiers: LR, ShallowDT, LoCapANN and HiCapANN. Thus, a high-data regime can tolerate duplicates by remaining more agnostic to the model biases.* Increasing the amount of duplication for a high data regime ($|r|_t=10k$) has a marginal impact on accuracy. *Thus, even high duplication has a marginal impact in the high-data regime.* We present the corresponding accuracy plots of the impact of duplicates with data regime changes on the other classifiers in Appendix.

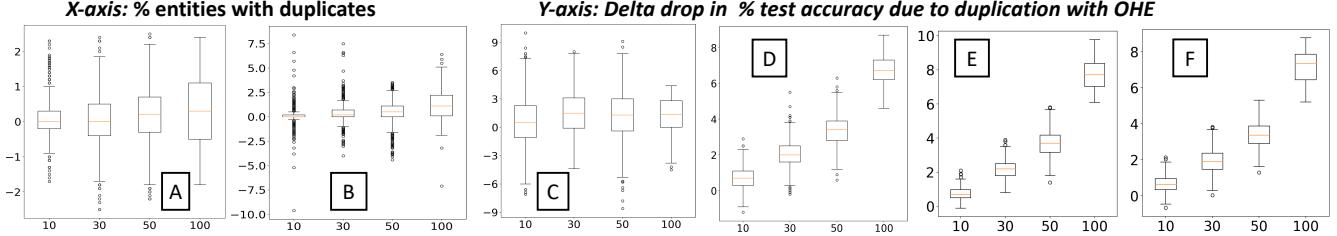


Figure 5: Results with OHE for (A) LR (B) ShallowDT (C) LoCapANN (D) HiCapANN (E) SVM (F) XGBoost with same setup as Figure 14(B).

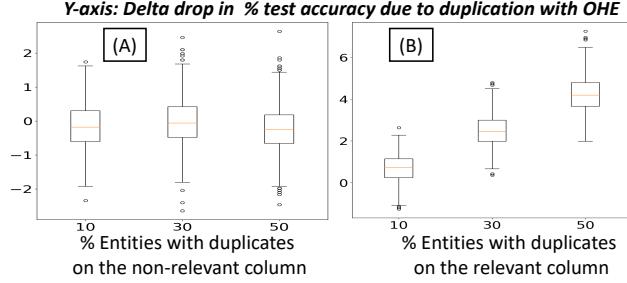


Figure 6: HiCapRF results. Vary $|ED|/|E|$, while fixing $(|\mathcal{A}|, |r|_t, \text{occ}(D_k), |D_k|) = (4, 5000, 25, 1)$. Duplicates introduced on the column with (A) non-positive Relevancy (noisy column) (B) high Relevancy (predictive column).

7.3.2 Varying parameters controlling the amount of duplication. Figure 14 (B-D) shows how different duplication parameters influence HiCapRF. We notice a clear trend: *the drop in accuracy with HiCapRF rises with the increase in any of the three duplication controlling parameters, $|ED|/|E|$, $\text{occ}(D_k)$, and $|D_k|$* . Among the three parameters, $|ED|/|E|$ has the most drastic effect on HiCapRF. The effects of the increase in $|D_k|$ are less pronounced because all other parameters including $\text{occ}(D_k)$ are kept fixed. Thus, there exist more duplicates for the same occurrence. *Interestingly, we find from Figure 14 that StrE is more robust to duplicates than OHE regardless of the parameter being varied, as the delta drop in accuracy with StrE is comparatively lower, although significant in high duplication cases.*

Figure 5 presents how a key EV affects other classifiers. We find that all high-bias models behave similarly as they show a marginal drop in accuracy even when all entities are diluted with duplicates. In contrast, HiCapANN exhibits similar behavior as HiCapRF when $|ED|/|E|$ is increased. Note that the absolute accuracies of the high-bias approaches are lower than that of high-capacity ones. *Overall, both high-capacity models are more susceptible to the adverse performance impact of duplicates than high-bias approaches.* We notice the same trend as other EVs, $\text{occ}(D_k)$ and $|D_k|$ are varied. We present the corresponding accuracy plots with them in Appendix.

7.3.3 Varying properties of duplicates being mapped to “Others.”

We study how duplicates that do not arise in the train set but are present in the test set (say, during deployment) can impact ML. We modify and repeat our duplication process on just the test set while keeping the train set intact. We introduce just one duplicate in the test set that gets mapped to “*Others*.” Figure 14 (E-F) presents the results on HiCapRF with OHE where $|ED|/|E|$ and $\text{occ}(D_k)$ are varied. *We find that the delta drop in accuracies with all parameters are even more higher than the corresponding delta drops when both train*

and test set were duplicated (Figure 14 (B-C)). This simply suggests that the presence of unwarranted duplicates during the test can cause downstream ML to suffer significantly.

7.3.4 Varying column Relevancy. We now study low vs. high *Relevancy* setting with a slight twist in our simulation. We introduce an additional noisy column in the clean dataset: All except one column participates in CPT. Thus, we have the presence of both high and low *Relevancy* columns. We introduce duplicates in both types of columns one at a time. Figure 6 present results. *We find that duplication on a highly relevant column has a significant adverse impact on HiCapRF performance. In contrast, the impact is negligible when duplicates are introduced over the noisy column. Even increasing the amount of duplication creates no impact with the low relevancy column. We observe the same trend with HiCapANN.*

7.4 Explanations and Takeaways

We now intuitively explain the behavior of ML classifiers in presence of duplicates. We check the generalization ability of ML models with the overfitting gap (as shown in Figure 7). We find that the delta accuracy drop (Figure 14) closely follows the increase in the overfitting gap due to duplicates with both high-capacity models, HiCapRF and HiCapANN. That is, the increase in overfitting or variance with duplicates explains the accuracy drop we see. *Thus, duplicates can negatively impact the generalization capability of high-capacity models, which are prone to overfitting.* However, as the number of training examples rises, the overfitting subsides. This explains our trends in the high-data regime.

We find that LR exhibits no amount of extra overfitting with duplicates. This is because the VC dimension of LR is linear in the number of features. As the dimensionality of the feature space expands with duplicates, VC dimension of LR expands. We get an expanded logistic hypothesis space with duplication that is a super-set of the true logistic hypothesis space. Thus, a larger hypothesis space can potentially lead to more variance unless the true concept is simple enough to recover in an expanded feature space. We check the weights of the hyperplane learned with LR in presence of duplicates where a higher weight indicates higher importance. We find that the absolute weights of duplicate features are often close to zero. *This suggests that the LR can learn the true concept by completely ignoring the extra dimensions. Thus, the variance does not rise.* HiCapRF with OHE makes many binary splits on the data to recover the true concept, causing the tree to fully grow to the restricted height. Chances of further overfitting with duplicates are reduced with a limited height. *This explains why a set-based split with StrE is more robust than binary splits with OHE as it allows to pack more category splits within the same tree height.*

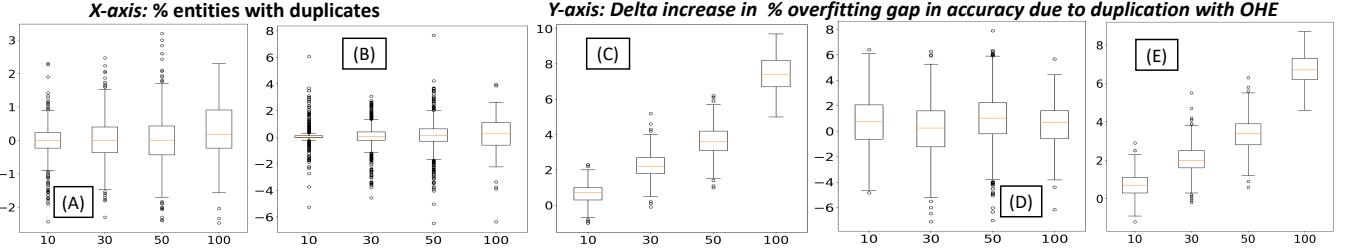


Figure 7: Simulation results on (A) LR (B) ShallowDT (C) HiCapRF (D) LoCapANN (E) HiCapANN (with the same setup as Figure 14(B)).

8 DISCUSSION

8.1 Public Release

We release a public repository on GitHub with our entire benchmark suite [4]. This includes our entire labeled dataset of 1262 *Categorical* columns along with entities in them annotated with corresponding duplicates and their raw CSV files. We also release the code to run downstream and benchmark suites.

8.2 Utility of our Labeled Data

Besides the utility of our labeled data for empirically benchmarking the impact of *Categorical* duplicates in Section 6, it also serve as a critical artifact to enable researchers in addressing many open questions. We highlight two important research directions below.

a. Design accurate methods for category deduplication. Although *Categorical* duplicates can often impact ML accuracy substantially, many existing open source AutoML tools such as AutoGluon [46] and TransmogrifAI [12] do not support an automated deduplication workflow. Cleaning duplicates manually or using ad hoc rules can be slow and error-prone for many users, especially non-technical lay users who were promised an end-to-end automation of the entire ML workflow. *Our labeled dataset will lead to an objective assessment of the accuracy of automation of different deduplication approaches.* Moreover, this will serve towards building supervised learning-based approach to automate the category deduplication task itself. In fact, one such hand-labeled data lead to highly accurate supervised ML approaches and even outperformed the existing industrial-strength tools for ML feature type inference [77].

b. Theoretical quantification. Our empirical study suggests that *Categorical* duplicates can increase variance since the hypothesis space of the model can grow. This opens up several research questions at the intersection of ML theory and data management: Is it possible to establish bounds on the increase in variance using VC-dimension theory [86]? Can we set up a decision rule to formally characterize when category deduplication would be needed? Our labeled data can be a key enabler to empirically validate the theory.

8.3 Takeaways for ML Practitioners

We find that the presence of *Categorical* duplicates can potentially impact downstream ML accuracy significantly. The amount of impact can be characterized by multiple EVs that interact in non-trivial ways. It is not always possible to disentangle the impact on ML with each EV individually. However, our empirical analyses can provide insights into when cleaning effort would be more or less beneficial. The current practice among ML practitioners

to handle *Categorical* duplicates is largely ad hoc rule-based and oblivious to many variables. While some of the presented insights have remained folklore for practitioners, our work presents the first systematic scientific study and put this on an empirically rigorous footing. We now give general guidelines and actionable insights to help them prioritise their category deduplication effort and also potentially design better end-to-end automation pipelines.

a. Make ML workflows less susceptible to the adverse performance impact of *Categorical* duplicates. LR is less prone to overfitting than the high-capacity models (RF, XGB, SVM, and ANN) when *Categorical* duplicates arise. This is because, as duplicates increase feature dimensionality of *Categoricals*, LR can completely ignore the extra dimensions of duplicates by setting their weights close to 0, making them overfit less. Also, *StrE* is relatively more robust than *OHE* when using RF. Moreover, *SimE*, *TransE*, and *TABBIE* are significantly more robust from *Categorical* duplicates compared to *OHE* and *StrE*. Tabular representation learning method, *TABBIE* doesn't offer significant benefit in terms of robustness to duplicates over the pre-trained embedding method *TransE*. Moreover, unseen *Categorical* duplicates that arise during the deployment phase can degrade ML performance with *OHE* or *StrE*. Overall, *Similarity* encoding and Logistic Regression or Transformer/TABBIE embeddings can be utilized by ML practitioners and AutoML developers if they desire to guard their pipelines against any adverse drop in ML performance from likely *Categorical* duplicates. Moreover, the impact of duplicates get mitigated in a higher-data regime compared to a low-data regime. Thus, whenever possible, one can consider getting more train data to offset their impact by trading off runtime.

ML practitioners who are considering to use category deduplication approaches can prioritize their efforts cleaning only specific kinds of duplicates: Abbreviations, Synonyms, and presence of Extra Information in duplicates for *SimE* with a high-capacity model and, Misspellings and Abbreviations types for *TABBIE*. The behavior of ML models with *OHE* and *StrE* do not vary considerably with a specific duplication type. In addition, *TransE* manifests significant robustness across all duplication types which makes it a reliable encoding choice for ML practitioners.

b. Track the overfitting gap of ML models. Category deduplication can reduce the overfitting caused by *Categorical* duplicates on ML. Thus, cleaning *Categorical* duplicates may not be worthwhile if the overfitting gap is already low on the raw data. Monitoring and presenting it as an auxiliary metric to the AutoML user can provide them with more confidence about the downstream performance.

REFERENCES

- [1] Accessed Dec 5, 2023. Data at <https://www.kaggle.com/datasets/jboysen/us-perm-visas-and-sourced-from> and sourced from <https://www.dol.gov/agencies/eta/foreign-labor/programs/permanent>.
- [2] Accessed Dec 5, 2023. Data at <https://www.kaggle.com/datasets/puppygogo/sa-dog-ownership-sample> and sourced from data.gov.au.
- [3] Accessed July 1, 2023. FlairNLP: Framework for state-of-the-art NLP, <https://github.com/flairnlp/flair>.
- [4] Accessed July 1, 2023. Github Repository for studying the impact of Cleaning Category Duplicates on ML, <https://github.com/pvn25/CategDupsRepo>.
- [5] Accessed July 1, 2023. Google AutoML Tables Cleaning Duplicates User Guidelines. https://cloud.google.com/automl-tables/docs/data-best-practices#make_sure_your_categorical_features_are_accurate_and_clean
- [6] Accessed July 1, 2023. Google AutoML Tables Data Prep User Guidelines, <https://cloud.google.com/automl-tables/docs/data-best-practices>.
- [7] Accessed July 1, 2023. Google Cloud AutoML, <https://cloud.google.com/automl>.
- [8] Accessed July 1, 2023. H2O Driverless AI, <https://www.h2o.ai/products/h2o-driverless-ai>.
- [9] Accessed July 1, 2023. H2o.AI, <https://www.h2o.ai>.
- [10] Accessed July 1, 2023. Microsoft AutoML, <https://azure.microsoft.com/en-us/services/machine-learning/automatedml/>.
- [11] Accessed July 1, 2023. Similarity Encoder Library, https://github.com/dirty-cat/dirty_cat.
- [12] Accessed July 1, 2023. TransmogrifAI: Automated Machine Learning for Structured Data, <https://transmogrif.ai>.
- [13] Accessed July 1, 2023. Trifacta: Data Wrangling Tools & Software, <https://www.trifacta.com/>.
- [14] Accessed July 1, 2023. <https://data.ca.gov/>.
- [15] Accessed July 1, 2023. <https://data.ca.gov/dataset/tsm-habitat-rapid-assessment-survey-2016-ds28271>.
- [16] Accessed July 1, 2023. <https://datacatalog.hsls.pitt.edu/>.
- [17] Accessed July 1, 2023. <https://datacatalog.hsls.pitt.edu/dataset/77>.
- [18] Accessed July 1, 2023. <https://data.cityofchicago.org/>.
- [19] Accessed July 1, 2023. <https://data.cityofchicago.org/Buildings/Vacant-and-Abandoned-Buildings-Violations/kc9i-wq85>.
- [20] Accessed July 1, 2023. <https://data.cityofchicago.org/Transportation/Relocated-Vehicles/5k2z-suxx>.
- [21] Accessed July 1, 2023. <https://data.ny.gov/>.
- [22] Accessed July 1, 2023. <https://data.ny.gov/Energy-Environment/Utility-Company-Customer-Service-Response-Index-CS-w3b5-8aqf>.
- [23] Accessed July 1, 2023. <https://everydaydata.co/>.
- [24] Accessed July 1, 2023. <https://everydaydata.co/2017/02/07/hacker-news-part-one.html>.
- [25] Accessed July 1, 2023. [https://github.com/fivethirtyeight/data/](https://github.com/fivethirtyeight/data).
- [26] Accessed July 1, 2023. <https://github.com/fivethirtyeight/data/tree/master/region-survey>.
- [27] Accessed July 1, 2023. [https://maxcandocia.com/article/2018/Oct/22/trick-or-treating-ages/](https://maxcandocia.com/article/2018/Oct/22/trick-or-treating-ages).
- [28] Accessed July 1, 2023. <https://osmihelp.org/>.
- [29] Accessed July 1, 2023. <https://osmihelp.org/research>.
- [30] Accessed July 1, 2023. <https://transparentcalifornia.com/salaries/san-francisco/>.
- [31] Accessed July 1, 2023. <https://www.asdcode.de/2021/01/it-salary-survey-december-2020.html>.
- [32] Accessed July 1, 2023. <https://www.kaggle.com/definitelyliliput/rawscores>.
- [33] Accessed July 1, 2023. <https://www.kaggle.com/milomuscio/wifi-study>.
- [34] Accessed July 1, 2023. <https://www.kaggle.com/pushpaltayal/etailing-customer-survey-in-india>.
- [35] Richard C. Angell, George E. Freund, and Peter Willett. 1983. Automatic Spelling Correction Using a Trigram Similarity Measure. *Inf. Process. Manag.* 19, 4 (1983), 255–261. [https://doi.org/10.1016/0306-4573\(83\)90022-5](https://doi.org/10.1016/0306-4573(83)90022-5)
- [36] Adel Ardalan, Derek Paulsen, Amanpreet Singh Saini, Walter Cai, and AnHai Doan. 2021. Toward Data Cleaning with a Target Accuracy: A Case Study for Value Normalization. *CoRR* abs/2101.05308 (2021). arXiv:2101.05308 <https://arxiv.org/abs/2101.05308>
- [37] Peter L Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian. 2019. Nearly-tight VC-dimension and pseudodimension bounds for piecewise linear neural networks. *The Journal of Machine Learning Research* 20, 1 (2019), 2285–2301.
- [38] Paul Suganthan G. C., Adel Ardalan, AnHai Doan, and Aditya Akella. 2018. Smurf: Self-Service String Matching Using Random Forests. *Proc. VLDB Endow.* 12, 3 (2018), 278–291. <https://doi.org/10.14778/3291264.3291272>
- [39] Patricio Cerdá, Gaël Varoquaux, and Balázs Kégl. 2018. Similarity encoding for learning with dirty categorical variables. *Machine Learning* 107, 8 (2018), 1477–1494.
- [40] Peter Christen. 2012. *Data Matching - Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection*. Springer. <https://doi.org/10.1007/978-3-642-31164-2>
- [41] Xu Chu, Ihab F. Ilyas, Sanjay Krishnan, and Jiannan Wang. 2016. Data Cleaning: Overview and Emerging Challenges. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 – July 01, 2016*, Fatma Özcan, Georgia Koutrika, and Sam Madden (Eds.). ACM, 2201–2206. <https://doi.org/10.1145/2882903.2912574>
- [42] Anamaria Crisan and Brittany Fiore-Gartland. 2021. Fits and Starts: Enterprise Use of AutoML and the Role of Humans in the Loop. In *CHI ’21: CHI Conference on Human Factors in Computing Systems, Virtual Event / Yokohama, Japan, May 8–13, 2021*, Yoshifumi Kitamura, Aaron Quigley, Katherine Isbister, Takeo Igashira, Pernille Bjørn, and Steven Mark Drucker (Eds.). ACM, 601:1–601:15. <https://doi.org/10.1145/3411764.3445775>
- [43] Sanjib Das, AnHai Doan, Paul Suganthan G. C., Chaitanya Gokhale, Pradap Konda, Yash Govind, and Derek Paulsen. Accessed July 1, 2023. The Magellan Data Repository. <https://sites.google.com/site/anhaidgroup/useful-stuff/the-magellan-data-repository?authuser=0>.
- [44] Sanjib Das, Paul Suganthan GC, AnHai Doan, Jeffrey F Naughton, Ganesh Krishnan, Rohit Deep, Esteban Arcuate, Vijay Raghavendra, and Youngchoon Park. 2017. Falcon: Scaling up hands-off crowdsourced entity matching to build cloud services. In *Proceedings of the 2017 ACM International Conference on Management of Data*. 1431–1446.
- [45] Dong Deng, Wenbo Tao, Ziwasch Abedjan, Ahmed K. Elmagarmid, Ihab F. Ilyas, Guoliang Li, Samuel Madden, Mourad Ouzzani, Michael Stonebraker, and Nan Tang. 2019. Unsupervised String Transformation Learning for Entity Consolidation. In *35th IEEE International Conference on Data Engineering, ICDE 2019, Macao, China, April 8–11, 2019*. IEEE, 196–207. <https://doi.org/10.1109/ICDE.2019.00026>
- [46] Nick Erickson, Jonas Mueller, Alexander Shirkov, Hang Zhang, Pedro Larroy, Mu Li, and Alexander J. Smola. 2020. AutoGluon-Tabular: Robust and Accurate AutoML for Structured Data. *CoRR* abs/2003.06505 (2020). arXiv:2003.06505 <https://arxiv.org/abs/2003.06505>
- [47] Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Tobias Springenberg, Manuel Blum, and Frank Hutter. 2015. Efficient and Robust Automated Machine Learning. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7–12, 2015, Montreal, Quebec, Canada*, Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett (Eds.), 2962–2970. <http://papers.nips.cc/paper/5872-efficient-and-robust-automated-machine-learning>
- [48] Pieter Gijsbers, Erin LeDell, Janek Thomas, Sébastien Poirier, Bernd Bischl, and Joaquin Vanschoren. 2019. An Open Source AutoML Benchmark. *arXiv preprint arXiv:1907.00909* (2019).
- [49] Chaitanya Gokhale, Sanjib Das, AnHai Doan, Jeffrey F Naughton, Narasimhan Rampalli, Jude Shavlik, and Xiaojin Zhu. 2014. Corleone: Hands-off crowdsourcing for entity matching. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*. 601–612.
- [50] Margherita Grandini, Enrico Bagli, and Giorgio Visani. 2020. Metrics for multi-class classification: an overview. *arXiv preprint arXiv:2008.05756*.
- [51] Trevor Hastie, Jerome H. Friedman, and Robert Tibshirani. 2001. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer. <https://doi.org/10.1007/978-0-387-21606-5>
- [52] Yeye He, Xu Chu, Kris Ganjam, Yudian Zheng, Vivek Narasayya, and Surajit Chaudhuri. 2018. Transform-Data-by-Example (TDE): An Extensible Search Engine for Data Transformations. *Proceedings of the VLDB Endowment* 11, 10 (2018), 1165–1177.
- [53] Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren (Eds.). 2019. *Automated Machine Learning - Methods, Systems, Challenges*. Springer. <https://doi.org/10.1007/978-3-030-05318-5>
- [54] Nick Hynes, D Sculley, and Michael Terry. 2017. The Data Linter: Lightweight, Automated Sanity Checking for ML Data Sets. In *NIPS MLSys Workshop*.
- [55] Hiroshi Iida, Dung Thai, Varun Manjunatha, and Mohit Iyer. 2021. TABIE: Pretrained Representations of Tabular Data. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 3446–3456.
- [56] Zhongjun Jin, Michael R Anderson, Michael Cafarella, and Hosagrahara V Jagadish. 2017. Foofah: A Programming-By-Example System for Synthesizing Data Transformation Programs. In *Proceedings of the 2017 ACM International Conference on Management of Data*. 1607–1610.
- [57] Mayank Kejriwal and Daniel P. Miranker. 2015. Semi-supervised Instance Matching Using Boosted Classifiers. In *The Semantic Web: Latest Advances and New Domains - 12th European Semantic Web Conference, ESWC 2015, Portoroz, Slovenia, May 31 – June 4, 2015. Proceedings (Lecture Notes in Computer Science)*, Fabien Gandon, Marta Sabou, Harald Sack, Claudia d’Amato, Philippe Cudré-Mauroux, and Antoine Zimmermann (Eds.), Vol. 9088. Springer, 388–402. https://doi.org/10.1007/978-3-319-18818-2_24
- [58] Pradap Venkatraman Konda. 2018. *Magellan: Toward building entity matching management systems*. The University of Wisconsin-Madison.
- [59] Sanjay Krishnan, Michael J. Franklin, Ken Goldberg, Jiannan Wang, and Eugene Wu. 2016. ActiveClean: An Interactive Data Cleaning Framework For Modern Machine Learning. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 – July 01, 2016*, Fatma Özcan, Georgia Koutrika, and Sam Madden (Eds.). ACM, 2117–2120.

- <https://doi.org/10.1145/2882903.2899409>
- [60] Sanjay Krishnan, Michael J. Franklin, Ken Goldberg, and Eugene Wu. 2017. BoostClean: Automated Error Detection and Repair for Machine Learning. *CoRR abs/1711.01299* (2017). arXiv:1711.01299 <http://arxiv.org/abs/1711.01299>
- [61] Peng Li, Xiang Cheng, Xu Chu, Yeye He, and Surajit Chaudhuri. 2021. Auto-FuzzyJoin: Auto-Program Fuzzy Similarity Joins Without Labeled Examples. In *SIGMOD '21: International Conference on Management of Data, Virtual Event, China, June 20–25, 2021*, Guoliang Li, Zhanhuai Li, Stratos Idreos, and Divesh Srivastava (Eds.). ACM, 1064–1076. <https://doi.org/10.1145/3448016.3452824>
- [62] Peng Li, Xi Rao, Jennifer Blase, Yue Zhang, Xu Chu, and Ce Zhang. 2021. CleanML: A Study for Evaluating the Impact of Data Cleaning on ML Classification Tasks. In *37th IEEE International Conference on Data Engineering, ICDE 2021, Chania, Greece, April 19–22, 2021*. IEEE, 13–24. <https://doi.org/10.1109/ICDE51399.2021.00009>
- [63] Yuliang Li, Jinfeng Li, Yoshihiko Suhara, AnHai Doan, and Wang-Chiew Tan. 2020. Deep Entity Matching with Pre-Trained Language Models. *Proc. VLDB Endow.* 14, 1 (2020), 50–60. <https://doi.org/10.14778/3421424.3421431>
- [64] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [65] David Maier. 1983. *The theory of relational databases*. Vol. 11. Computer science press Rockville.
- [66] Venkata Vamsikrishna Meduri, Lucian Popa, Prithviraj Sen, and Mohamed Sarwat. 2020. A Comprehensive Benchmark Framework for Active Learning Methods in Entity Matching. In *Proceedings of the 2020 International Conference on Management of Data, SIGMOD Conference 2020, online conference [Portland, OR, USA], June 14–19, 2020*, David Maier, Rachel Pottinger, AnHai Doan, Wang-Chiew Tan, Abdussalam Alawini, and Hung Q. Ngo (Eds.). ACM, 1133–1147. <https://doi.org/10.1145/3318464.3380597>
- [67] Sidharth Mudgal, Han Li, Theodoros Rekatsinas, AnHai Doan, Youngchoon Park, Ganesh Krishnan, Rohit Deep, Esteban Arcaute, and Vijay Raghavendra. 2018. Deep Learning for Entity Matching: A Design Space Exploration. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10–15, 2018*, Gautam Das, Christopher M. Jermaine, and Philip A. Bernstein (Eds.). ACM, 19–34. <https://doi.org/10.1145/3183713.3196926>
- [68] Avanika Narayan, Ines Chami, Laurel Orr, Simran Arora, and Christopher Ré. 2022. Can Foundation Models Wrangle Your Data? *arXiv preprint arXiv:2205.09911* (2022).
- [69] Felix Neutatz, Binger Chen, Ziawasch Abedjan, and Eugene Wu. 2021. From Cleaning before ML to Cleaning for ML. *IEEE Data Eng. Bull.* 44, 1 (2021), 24–41. <http://sites.computer.org/debull/A21mar/p24.pdf>
- [70] Felix Neutatz, Binger Chen, Yazan Alkhateb, Jingwen Ye, and Ziawasch Abedjan. 2022. Data Cleaning and AutoML: Would an optimizer choose to clean? *Datenbank-Spektrum* (2022), 1–10.
- [71] Randal S. Olson, Nathan Bartley, Ryan J. Urbanowicz, and Jason H. Moore. 2016. Evaluation of a Tree-based Pipeline Optimization Tool for Automating Data Science. In *Proceedings of the 2016 on Genetic and Evolutionary Computation Conference, Denver, CO, USA, July 20–24, 2016*, Tobias Friedrich, Frank Neumann, and Andrew M. Sutton (Eds.). ACM, 485–492. <https://doi.org/10.1145/2908812.2908918>
- [72] Fatemeh Panahi, Wentao Wu, AnHai Doan, and Jeffrey F. Naughton. 2017. Towards Interactive Debugging of Rule-based Entity Matching. In *Proceedings of the 20th International Conference on Extending Database Technology, EDBT 2017, Venice, Italy, March 21–24, 2017*, Volker Markl, Salvatore Orlando, Bernhard Mitschang, Periklis Andritsos, Kai-Uwe Sattler, and Sebastian Breß (Eds.). OpenProceedings.org, 354–365. <https://doi.org/10.5441/002/edbt.2017.32>
- [73] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research* 12 (2011), 2825–2830.
- [74] Jinglin Peng, Weiyuan Wu, Brandon Lockhart, Song Bian, Jing Nathan Yan, Linghao Xu, Zhixuan Chi, Jeffrey M. Rzeszotarski, and Jiannan Wang. 2021. DataPrep.EDA: Task-Centric Exploratory Data Analysis for Statistical Modeling in Python. In *Proceedings of the 2021 International Conference on Management of Data (SIGMOD '21), June 20–25, 2021, Virtual Event, China*. Data (SIGMOD '21), June 20–25, 2021, Virtual Event, China.
- [75] J. Ross Quinlan. 1986. Induction of decision trees. *Machine learning* 1, 1 (1986), 81–106.
- [76] Theodoros Rekatsinas, Xu Chu, Ihab F. Ilyas, and Christopher Ré. 2017. Holoclean: Holistic Data Repairs with Probabilistic Inference. *Proc. VLDB Endow.* 10, 11 (2017), 1190–1201. <https://doi.org/10.14778/3137628.3137631>
- [77] Vraj Shah, Jonathan Lacanlale, Premanand Kumar, Kevin Yang, and Arun Kumar. 2021. Towards Benchmarking Feature Type Inference for AutoML Platforms. In *Proceedings of the 2021 International Conference on Management of Data*. 1584–1596.
- [78] Rohit Singh, Venkata Vamsikrishna Meduri, Ahmed K. Elmagarmid, Samuel Madden, Paolo Papotti, Jorge-Arnulfo Quiané-Ruiz, Armando Solar-Lezama, and Nan Tang. 2017. Generating Concise Entity Matching Rules. In *Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD Conference 2017, Chicago, IL, USA, May 14–19, 2017*, Semih Salihoglu, Wenchao Zhou, Rada Chirkova, Jun Yang, and Dan Suciu (Eds.). ACM, 1635–1638. <https://doi.org/10.1145/3035918.3058739>
- [79] Rohit Singh, Venkata Vamsikrishna Meduri, Ahmed K. Elmagarmid, Samuel Madden, Paolo Papotti, Jorge-Arnulfo Quiané-Ruiz, Armando Solar-Lezama, and Nan Tang. 2017. Synthesizing Entity Matching Rules by Examples. *Proc. VLDB Endow.* 11, 2 (2017), 189–202. <https://doi.org/10.14778/3149193.3149199>
- [80] Survey. Accessed July 1, 2023. 2021 State of Data Science and Machine Learning. <https://www.kaggle.com/kaggle-survey-2021>.
- [81] Survey. Accessed July 1, 2023. 2022 State of Data Science and Machine Learning Results Visualized. <https://www.kaggle.com/paultimothymooney/kaggle-survey-2022-all-results>.
- [82] Survey. Accessed July 1, 2023. 2022 State of Data Science and Machine Learning. <https://www.kaggle.com/kaggle-survey-2022>.
- [83] Nan Tang, Ju Fan, Fangyi Li, Jianhong Tu, Xiaoyong Du, Guoliang Li, Samuel Madden, and Mourad Ouzzani. 2021. RPT: Relational Pre-trained Transformer Is Almost All You Need towards Democratizing Data Preparation. *Proc. VLDB Endow.* 14, 8 (2021), 1254–1261. <https://doi.org/10.14778/3457390.3457391>
- [84] Chris Thornton, Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. 2013. Auto-WEKA: Combined Selection and Hyperparameter Optimization of Classification Algorithms. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 847–855.
- [85] Jianhong Tu, Ju Fan, Nan Tang, Peng Wang, Guoliang Li, Xiaoyong Du, Xiaofeng Jia, and Song Gao. 2023. Unicorn: A unified multi-tasking model for supporting matching tasks in data integration. *Proceedings of the ACM on Management of Data* 1, 1 (2023), 1–26.
- [86] Vladimir Naumovich Vapnik. 2000. *The Nature of Statistical Learning Theory, Second Edition*. Springer.
- [87] David Vos, Till Döhmen, and Sebastian Schelter. 2022. Towards Parameter-Efficient Automation of Data Wrangling Tasks with Prefix-Tuning. In *NeurIPS 2022 First Table Representation Workshop*.
- [88] Renzhi Wu, Sanya Chaba, Saurabh Sawlani, Xu Chu, and Saravanan Thirumuruganathan. 2020. ZeroER: Entity Resolution using Zero Labeled Examples. In *Proceedings of the 2020 International Conference on Management of Data, SIGMOD Conference 2020, online conference [Portland, OR, USA], June 14–19, 2020*, David Maier, Rachel Pottinger, AnHai Doan, Wang-Chiew Tan, Abdussalam Alawini, and Hung Q. Ngo (Eds.). ACM, 1149–1164. <https://doi.org/10.1145/3318464.3389743>
- [89] Doris Xin, Eva Yiwei Wu, Doris Jung Lin Lee, Niloufar Salehi, and Aditya G. Parameswaran. 2021. Whither AutoML? Understanding the Role of Automation in Machine Learning Workflows. In *CHI '21: CHI Conference on Human Factors in Computing Systems, Virtual Event / Yokohama, Japan, May 8–13, 2021*, Yoshifumi Kitamura, Aaron Quigley, Katherine Isbister, Takeo Igarashi, Pernille Bjørn, and Steven Mark Drucker (Eds.). ACM, 83:1–83:16. <https://doi.org/10.1145/3411764.3445306>
- [90] Chen Zhao and Yeye He. 2019. Auto-EM: End-to-end Fuzzy Entity-Matching using Pre-trained Deep Models and Transfer Learning. In *The World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13–17, 2019*, Ling Liu, Ryen W. White, Amin Mantrach, Fabrizio Silvestri, Julian J. McAuley, Ricardo Baeza-Yates, and Leila Zia (Eds.). ACM, 2413–2424. <https://doi.org/10.1145/3308558.3313578>

APPENDIX

9 HAND-LABELED DATA

9.1 Existing Labeled Datasets

Many works have introduced labeled data for related tasks such as EM [43, 44, 58], SM [38, 43, 61], and resolving column-level inconsistencies [62]. Table 7 gives examples of pairs of duplicates from three different datasets in Magellan Data Repository [43], which is used in several prior works [38, 44, 49, 58]. Open domain attributes such as person names and addresses are not *Categorical* features for ML. Instead, such features are context-specific where either custom features are extracted or are completely dropped as they may not generalize for ML. Note that *Categorical* feature assumes mutually exclusive values from a known finite domain. Moreover, *Title* in *Citations* has rich semantic information and is typically used as a *Text* type feature. We find that a large fraction of the datasets used in prior works involve duplication in non-*Categorical* features. Thus, they are not relevant for us to study category deduplication. Although incidental *Categorical* duplicates do arise in a few datasets in [38, 62], we posit that we need a systematic benchmark to characterize and understand their impact on ML accuracy that prior works do not focus on. We focus exclusively on the *Categorical* features and curate the first large labeled dataset of entities annotated with duplicates within a *Categorical* column.

Table 7: Examples of labeled duplicates from Magellan Data Repository [43] with dataset and column names.

	Address	Phone number	Name
A. Restaurants	1929 Hillhurst Ave, Los Angeles, CA	(323) 644-0100	Alcove Cafe & Bakery
	1929 Hillhurst Ave, Los Angeles, CA 90027	(323) 644-0100	Alcove Cafe & Bakery
B. Citations	Author	Entry Type	Title
	David A. Cohn and Michael I. Jordan	article	Active Learning with Statistical Models
C. Researchers	Cohn, David A and Jordan, Michael I	article	Active learning with statistical models
	Name	X	
	alicia n aarnio		
	alicia nicole aarnio		

9.2 Data Statistics and Takeaways

There exist three parameters that quantify the amount of duplication within a column. (1) Fraction of entities that have at least one duplicate ($|ED|/|E|$). (2) Duplicate set size for all entities of the column (set D). (3) Duplicate occurrences $occ(\{D_{ki}\})$, $k \in [1, |ED|]$, $i \in [1, |D|]$. Figure 8 presents the cumulative distribution function (CDF) of these parameters over all columns in our labeled dataset that has at least 1 duplicate. The Y-axis in the CDF plots show the probability that the presented parameter takes a value less than or equal to X -axis values. We also report CDF of the % reduction in domain size of the columns with deduplication. We now summarize the key takeaways below.

(1) We first explain the worst-case scenario that can arise due to duplicates. We find that almost 17% of the columns that have duplicates have them in all of their entities! Furthermore, 7% of duplicates

across all columns occur at 50%, i.e., the representation of the duplicate and the true entity are same. Additionally, 1% of all entities have more than five duplicates. However, the presence of more than 10 duplicates per entity is quite unlikely (less than 0.5%). Finally, deduplication can reduce the domain size of the *Categorical* column substantially by up to 99%. Overall, we find that whenever duplicates arise in the column, they can occur quite often.

(2) We now discuss the presence of duplicates with the average case scenario. We find that whenever an entity is diluted with duplicates, almost 90% of the time they have one or two duplicates! Duplicate set sizes follow a long-tail distribution, most entities have small duplicate set sizes and very few entities have a lot of duplicates. This can make catching duplicates and deduplicating them particularly challenging, as they can go unnoticed. Moreover, the occurrence of duplicates approximately follows a uniform distribution, i.e., all occurrence values up to 50% are roughly equally likely. Finally, $|ED|/|E|$ values of 10-35% fall close to the median.

(3) We present how different duplication types (from Table 3 of the paper) are represented in our labeled data in Figure 9. We find that the *Difference of Special Characters* and *Capitalization* issues are the most common, while *Synonyms* and *Grammar* issues are less common. Moreover, whenever duplicates exist, 17% of the time they belong to more than one duplication type (maximum observed is 4 duplication types). Also, 19% of entities have duplicates that can be mapped to multiple types (maximum observed is 6 duplication types).

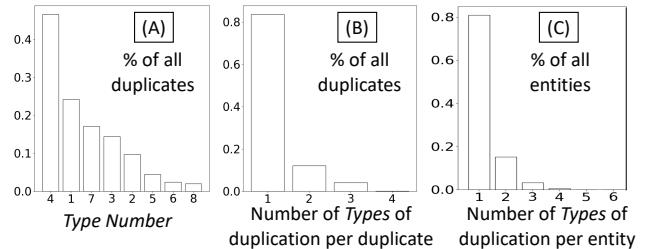


Figure 9: Histogram plots to illustrate how duplication types (from Table 3 of the paper) arise across all columns in all files. x and y denote x-axis and y-axis respectively. (A) $y\%$ duplicates have duplication of at least one x Type Number. (B) $y\%$ duplicates have x different duplication types within. (C) $y\%$ entities have duplicates with x different duplication types within.

10 DOWNSTREAM BENCHMARK

10.1 Dataset Statistics

Table 8 presents the statistics with different EVs that can potentially impact the ML performance over all 16 downstream datasets.

10.2 Methodology and Experimental Setup

We use scikit-learn[73], H2o[9], SimilarityEncoder[11], and FlairNLP[3] packages to employ *OHE*, *StrE*, *SimE*, and *TransformerE* resp. We use a standard grid search for hyper-parameter tuning of the models, with the grids described in detail below.

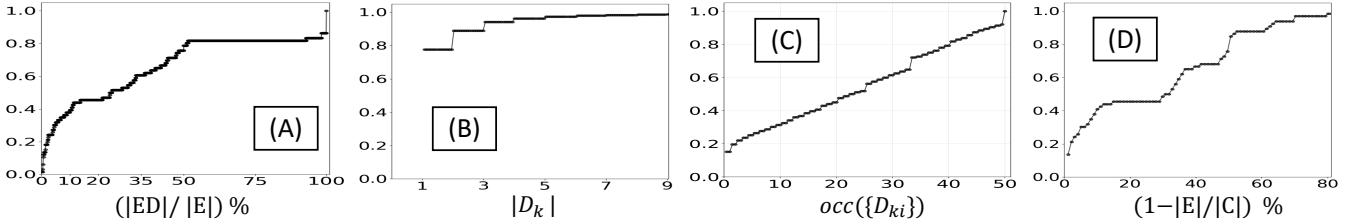


Figure 8: Cumulative distribution function (CDF) over all *Categorical* columns with at least one duplicate on (A) % entities that have at least one duplicate. (B) Duplicate set sizes over all $k \in [1, |ED|]$. The maximum duplicate set size is 148. (C) Duplicate set occurrences over all $k \in [1, |ED|], i \in [1, |D_k|]$. (D) % reduction in domain size with category deduplication.

Table 8: Statistics of the column containing duplicates in our downstream benchmark datasets. $|r|$, $|\mathcal{A}|$, and $|Y|$ are the total number of examples, number of columns, and number of target classes in the given dataset respectively. $|rC|$ denotes the number of training examples (averaged over 5 folds) per category of the set C . P is the fraction of $|E|$ (averaged across 5-folds) that has at least 1 duplicate being mapped to “*Others*” category in the validation set with *OHE* and *StrE*. We use colors green, blue, red with hand-picked thresholds to visually present and better interpret the cases where the amount of duplication is low ($1 - |E|/|C| < 0.25$), moderate ($1 - |E|/|C| > 0.25 \& < 0.50$), and high ($1 - |E|/|C| > 0.50$) respectively. We use the following thresholds with the same colors to better interpret the data regime: low ($|rC| < 5$), moderate ($|rC| > 5 \& < 25$), and high ($|rC| > 25$). Note that the data regime moves up with deduplication as category set size has shrunk.

Datasets	$ r $	$ \mathcal{A} $	$ Y $	Amount of Duplication					Data Regime		$P\%$
				$\frac{ ED }{ E }\%$	median $ D_k $	median $\text{occ}(\{D_{ki}\})$	$ C $	$1 - \frac{ E }{ C }\%$	$ rC $	$ rC $ after dedup (Increase w.r.t Raw)	
Midwest Survey	2778	29	9	33.1	2	4	1008	64	2.5	6.5 (2.6x)	23.6
Mental Health	1260	27	5	40	3.5	2.3	49	69.4	23.2	81.2 (3.5x)	25.3
Relocated Vehicles	3263	20	4	33.2	1	20	1097	35.8	2.5	3.8 (1.5x)	14.9
Health Sciences	238	101	4	36.4	2	6	56	60.7	3.6	8.3 (2.3x)	26.4
Salaries	1655	18	8	24	1	25	647	29.2	1.8	2.2 (1.2x)	10.9
TSM Habitat	2823	48	19	11	1	25	912	11.4	2.6	2.9 (1.1x)	14.6
EU IT	1253	23	5	24	1.5	12.5	256	34.8	3.9	5.9 (1.5x)	19.5
Halloween	292	55	6	31.3	2	11.1	163	50.9	1.5	3 (2x)	22.8
Utility	4574	13	95	38.4	1	20	199	30.7	16.2	24.3 (1.5x)	6.2
Mid or Feed	1006	78	5	21.4	6	0.8	37	62.2	20.6	59.7 (2.9x)	24.3
Wifi	98	9	2	30.3	2.5	12.5	69	52.2	1.3	2.5 (1.9x)	26.1
Etailing	439	44	5	47.8	4	5.9	71	67.6	5.3	14.3 (2.7x)	28.7
San Francisco	148654	13	2	10.7	1	25	2159	9.8	46.3	50.9 (1.1x)	3.2
Building Violations	22012	17	6	51	2	4.8	270	63	53.7	145 (2.7x)	4.4
US Labor	210287	25	4	25	2	20	1169	47	31	60.4 (2x)	24.5
Pet Registration	82545	14	2	37	2	25	789	44	58.2	105 (1.8x)	18.5

Logistic Regression: There is only one regularization parameter to tune: C . Larger the value of C , lower is the regularization strength, hence increasing the complexity of the model. The grid for C is set as $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 100, 10^3\}$.

Random Forest: There are two hyper-parameters to tune: *NumEstimator* and *MaxDepth*. *NumEstimator* is the number of trees in the forest. *MaxDepth* is the maximum depth of the tree. The grid

is set as follows: $\text{NumEstimator} \in \{5, 25, 50, 75, 100\}$ and $\text{MaxDepth} \in \{5, 10, 25, 50, 100\}$.

ANN: The artificial neural network (ANN) architecture comprises of 2 hidden units with 100 neurons each. We do L_2 regularization with the regularization parameter tuned using the following grid axis: $\{10^{-4}, 10^{-3}, 10^{-2}\}$.

SVM: We tune hyper-parameter C , which controls the cost of misclassification. The grid is set as $\{10^{-1}, 1, 10, 100, 10^3\}$.

Table 9: ML Accuracy Results using *Raw* with different ML models and encoding schemes. The accuracy lifts compared to *Raw* are shown in Table 5 of the paper.

Dataset	Random Forest			XGBoost			ANN			SVM		Logistic Regression		TABBIE + ANN-1L
	OHE	StrE	SimE	OHE	StrE	SimE	OHE	SimE	TransE	OHE	SimE	OHE	SimE	
Midwest Survey	49.1	59.2	64.9	54.4	55.1	64.7	54.7	63.4	56.8	51.7	67.6	57.2	66.7	57
Mental Health	47.9	47.8	47.4	44.5	43.7	43.5	42.4	43.2	45.1	47.5	46.8	46.9	46.3	41.1
Relocated Vehicles	72.5	81.3	88.3	73.2	78.1	87.6	83.6	89.6	77	78.4	87.7	82.9	88.4	76.5
Health Sciences	53.3	61.8	60	45.3	48.4	45.3	55.1	56.4	57.3	52.4	54.7	58.7	60	56.4
Salaries	64.7	69.6	94.6	66.6	69.4	63.6	22	19.9	25	55.8	55.8	30.4	32.4	29.8
TSM Habitat	71.2	84.1	71.2	80	80.4	89.1	50.7	50.7	35.3	61.4	61.4	50.7	50.7	41.3
EU IT	41.2	43.6	47.8	48.1	48.7	49.9	13.4	6.8	9.9	33.3	33.3	29.1	29.1	7.7
Halloween	40	36.2	34.7	43	32.8	45.7	41.9	43	41.5	35.8	32.1	42.6	49.8	40
Utility	58.8	46.3	43.2	87.3	84.7	86.1	65.1	73.2	82.1	45.7	45.9	42.4	43	80.8
Mid or Feed	40.2	35.7	36.2	38	38.2	41.7	34	32.7	33.5	39.3	39.2	40.5	41.5	33.3
Wifi	60	57.9	50.5	57.9	66.3	56.8	52.6	48.4	61.1	64.2	63.2	64.2	58.9	45.6
Etailing	40	44.5	38.2	41.6	38.4	41.6	40.2	37.2	36.6	42.7	42.3	41.1	38.9	47.9
San Francisco	83.4	83.9	86	85.9	85.9	86.2	86	86.1	85.6	86	86.1	86	85.5	85.9
Building Violations	97.5	97.3	97.6	98.1	98	98.1	97.2	97.4	97.6	92.3	92.4	91.6	91.9	97
US Labor	74.1	75	75.9	75.9	74.5	76	68.8	73.3	79.9	74	75.1	65.7	68.8	79.9
Pet Registration	91.8	91.8	92.2	91.5	91	90.8	92.2	92	91	87.2	88.2	88.4	89	91.1

Table 10: Comparison of three representative downstream models in terms of Macro F1 score with different *Categorical* encoding schemes on *Raw* (column with duplicates) vs. *Truth* (deduplicated column) data. Results for *Truth* are shown relative to the *Raw* as delta lift/drop in % F1 score. Green, blue, and red colors denote cases where the *Truth* F1 score relative to *Raw* is significantly higher, comparable, and significantly lower (error tolerance of 1%) respectively.

Dataset	Logistic Regression				Random Forest						ANN			
	OHE		SimE		OHE		StrE		SimE		OHE		SimE	
	Raw	Truth	Raw	Truth	Raw	Truth	Raw	Truth	Raw	Truth	Raw	Truth	Raw	Truth
Midwest Survey	55.7	+10.3	65.4	+2.7	44.9	+12.6	56.5	+11.7	63.4	+5	54.3	+9.7	63.3	+3.7
Mental Health	42	-0.6	40.1	+0.8	40.3	+0	39.3	-1.3	38	+0.8	39.3	+2.7	41.1	+0.5
Relocated Vehicles	82.8	+4	88.4	+0.4	71.6	+3.5	81.3	+3.7	88.3	-0.1	83.6	+3.6	89.5	+0
Health Sciences	56.1	+0.7	57.4	+2.2	51.5	+3.3	59.1	-0.5	59.2	-3.7	54.7	+5.4	56.6	+1.8
Salaries	27.4	+1.3	30.5	-2	57.6	+2.1	64.5	+1.9	93.8	+0.4	15.9	+3.4	14.7	+8.7
TSM Habitat	34.1	+0	34.1	+0	68.5	+0	82.2	+1.6	85.7	+0.6	24.6	+4	19.1	+12.3
EU IT	16.1	+0	16.1	+0	33.6	+1.1	36.8	+0.2	43.1	+2.7	9.3	-2.2	4.2	+4.8
Halloween	37.1	+3.8	45.7	+0.5	34.2	+3	33.2	+1.7	31.1	-4.9	38.8	+4.1	40.9	-0.6
Utility	37	+0.1	38.5	+0.3	58.2	+1.5	44.9	+1.4	41.8	+1.7	65.2	+2	73.4	+2.2
Mid or Feed	37.2	+0.2	39.1	-3.6	35.2	+1.8	26.6	-0.3	26.1	+2.6	33	+1.9	31.2	+0.4
Wifi	54.9	+1.5	50.7	+8.2	52.7	+8.5	54.3	-4.5	50.2	+1.9	51.6	+2.3	48.3	+2.6
Etailing	37.2	-2.3	37.5	-0.1	33.3	+3	36.3	+1.4	32.9	+3.1	39.4	-3.1	36.7	+0
San Francisco	85.9	-0.1	85.6	-0.1	83	+0.3	83.3	+0.3	86.1	-0.1	86	+0.1	86	+0.1
Building Violations	89.4	+0	89.3	+0.1	97.5	-0.1	97	+0.1	97.4	+0.1	97.5	+0.1	97.2	+0.4

XGBoost: We tune two parameters "NumEstimator" and "MaxDepth" with grids {5, 25, 50, 100} and {5, 25, 50, 100} respectively. We use the default values provided by XGBoost library for tree booster parameters such as learning rate

10.3 Results with Additional Evaluation Metrics

We check if using additional evaluation metrics such as F1 score, precision, and recall alter any empirical observations or conclusions

in Section 6.4.1 of the paper. Note that the micro average of precision, recall, and F1-score is identical to the accuracy of multi-class classification [50]. Thus, we use the macro average of precision, recall, and F1-score [50] as evaluation metrics and rerun our downstream benchmark suite. We check if evaluating with macro F1 score alters the conclusion made with % diagonal classification accuracy as the metric in regard to the varied EV.

Overall, we find that none of the empirical conclusions made with diagonal accuracy in the paper change even with these additional evaluation metrics. Table 10 presents the comparison of downstream models with different *Categorical* encoding schemes in terms of macro F1 scores. Table 11 showcases the aggregate statistics over all evaluation metrics. We confirm the validity of all observations O1-O8 made with the downstream benchmark suite with the additional evaluation metrics. As an example, we still find that the delta increases in macro F1-score, precision, and recall with *Truth* are higher with Random Forest and artificial neural network (ANN) than Logistic Regression. Moreover, we again find that *Similarity* encoding is more robust than other encoding schemes to tolerate duplicates.

11 SIMULATION STUDY RESULTS

11.1 Scenario Alla

This represents a complex joint distribution where all features obtained from \mathcal{A} determine Y .

11.1.1 Varying the data regime and the parameters that control the amount of duplication. Figure 11 shows the delta drop in classification accuracy with duplication relative to the ground truth clean dataset on all models with *OHE* as the number of training examples and the duplication variables are varied. Figure 12 shows the results as the EVs are varied for HiCapRF with *StrE*. We again note that a high-data regime is more robust to duplication than a low-data regime for both encoding schemes. All the high-bias approaches are more robust to duplication than the high-capacity models. Also, duplication variables can have significant adverse performance effects on high-capacity classifiers

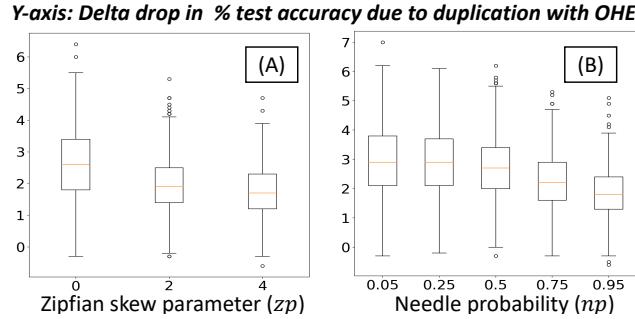


Figure 10: Effects of skew parameters on Alla simulation scenario for Random Forest with OHE. $|A|$ is preset to 3. (A) Vary Zipfian skew parameter zp of $|G_k|$, while fixing $(|r|_t, |ED|/|E|) = (3000, 30)$. (B) Vary needle probability parameter np of $occ(G_k)$, while fixing $(|r|_t, |ED|/|E|, zp) = (3000, 30, 2)$.

11.1.2 Skewness in the duplication parameters. Until now, we assumed that all entities in ED have identical duplicate set sizes $|D_k|$ and all duplicates in D_k are equally likely to occur. From our labeled data, we find that most entities have small duplicate set sizes and only a few entities have many duplicates. Also, some duplicates are more likely to occur than others in D_k . Thus, we relax these two assumptions and include distributions in $|D_k|$ and $occ(D_k)$ that can better represent the duplication process. We alter our duplication process and approximate $|D_k|$ with a long-tail Zipfian distribution and $occ(D_k)$ with a Needle-and-Thread distribution, varying the skew amount one at a time. *Overall, we find that none of our takeaways change or get invalidated with this setup.* We further explain below.

We now alter our duplication process to capture skewness in $|D_k|$ and $occ(D_k)$, varying one at a time. Figure 10 presents the results. We find that the delta drop in % accuracy due to duplicates remains significant regardless of the amount of skew in $|D_k|$ and $occ(D_k)$. With the Zipfian skew in $|D_k|$, the delta drop is highest at uniform distribution in $|D_k|$ (no skew setting) and marginally decreases as the skewness parameter is increased. On the other hand, when a needle-and-thread skew in $occ(D_k)$ is present, one duplicate from set D_k has a probability mass np (needle parameter). The remaining $1-np$ probability mass is uniformly distributed over the rest of the duplicates in D_k . We find that the delta drop due to duplicates decreases while still remaining significant when one duplicate value is more likely to occur than the rest (as np is increased). *Overall, the overarching conclusion from this analysis is that none of our results or takeaways change or get invalidated with this new setup.*

11.2 Scenario Hyperplane

Data generating process. We set up distribution with a true hyperplane to separates the classes. (1) We define and fix the normal vector of the hyperplane with weights, $W_i, 1 <= i <= |\mathcal{A}|$. Each weight W_i with cardinality $|E|$ is randomly sampled from a list of non-zero integers $([-5, 5] \setminus \{0\})$ without replacement. Note that the integer weights are chosen only to make the distance calculation simpler in step (3). The trends of our results do not change even if the weights are chosen from real number uniform distribution. (2) Construct $|r|$ tuples of \mathcal{A} by sampling values uniformly randomly from $|E|$. Thus, with fixed weights, the hyperplane over *One-hot* encoded example feature vectors is given by $\sum_{i=1}^{i=|\mathcal{A}|} W_i \cdot A_i = 0$. (3) Examples for which $\sum_{i=1}^{i=|\mathcal{A}|} W_i \cdot A_i >= 0$ are labeled positive ($Y=0$) and remaining examples are labeled negative ($Y=1$). This generates the true dataset where all columns have high *Relevancy*. We introduce duplicates in them by following the same duplication process as Section 7.2.

Results. Figure 13 shows the delta drop in accuracy due to duplicates with all models using *OHE*. We find that all high-bias approaches are again robust to the presence of duplicates even when all entities are diluted with duplicates. Interestingly, HiCapANN exhibits only marginal impact with duplicates. In contrast, duplicates affect HiCapRF significantly, especially in a high-duplication regime. We explain this interesting behavior in Section 7.4. We vary other variables such as other duplication parameters, the fraction of entities being mapped to “*Others*,” and column *Relevancy*. We confirm

Table 11: Summary stats over 16 downstream datasets in terms of Macro F1-score, precision, and recall to showcase the impact of deduplication on ML using different encodings.

F1 score	LR		Random Forest			ANN	
	OHE	SimE	OHE	StrE	SimE	OHE	SimE
% lift in accuracy with <i>Truth</i>							
Mean	1.4	0.7	2.9	1.4	0.7	2.4	2.6
Median	0.2	0.2	2	1.4	0.7	2.5	1.2
75 th percentile	1.5	0.7	3.2	1.9	2.5	3.9	3.4
Max	10.3	8.2	12.6	11.7	5	9.7	12.3
# downstream datasets where							
>1% lift wrt metric on <i>Truth</i>	5	3	10	8	6	10	7
Precision	LR		Random Forest			ANN	
	OHE	SimE	OHE	StrE	SimE	OHE	SimE
% lift in accuracy with <i>Truth</i>							
Mean	1.8	1.1	4.2	1.2	0.7	1.7	1.9
Median	0.9	0.3	3.2	1	0.5	1.9	1.7
75 th percentile	2.9	2.1	6.2	1.9	2.8	3.4	2.8
Max	10.3	8.1	14.7	11.9	5.8	9.8	9.8
# downstream datasets where							
>1% lift wrt metric on <i>Truth</i>	7	4	10	7	6	9	8
Recall	LR		Random Forest			ANN	
	OHE	SimE	OHE	StrE	SimE	OHE	SimE
% lift in accuracy with <i>Truth</i>							
Mean	1.6	0.9	2.4	1.9	0.7	2.3	2.7
Median	0.9	0.4	1.6	1.3	0.6	2.1	1.3
75 th percentile	1.6	1.1	2.4	2.1	2.7	4.1	3.7
Max	9.4	8.4	10.8	10	4.4	9.5	15
# downstream datasets where							
>1% lift wrt metric on <i>Truth</i>	7	4	10	9	6	9	7

the same trends that we saw with all models in AllA scenario, except with HiCapANN which behaves similar to LR than HiCapRF.

11.2.1 Varying the data regime and the parameters that control the amount of duplication. Figure 15 presents the delta drop in accuracy due to duplicates with all models using *OHE*. We again note that as the number of available training examples is increased, the delta drop in accuracy due to duplicates decreases for HiCapRF. Raising the other duplication parameters such as $|ED|/|E|$, $occ(D_k)$, $|D_k|$ also increases the adverse performance impact of duplicates on HiCapRF. Interestingly, we find that HiCapANN exhibits only a marginal amount of overfitting on the Hyperplane simulation scenario.

Thus, we do not see any impact due to duplicates on HiCapANN and also on all the high-bias approaches.

Y-axis: Delta drop in % test accuracy due to duplication with OHE

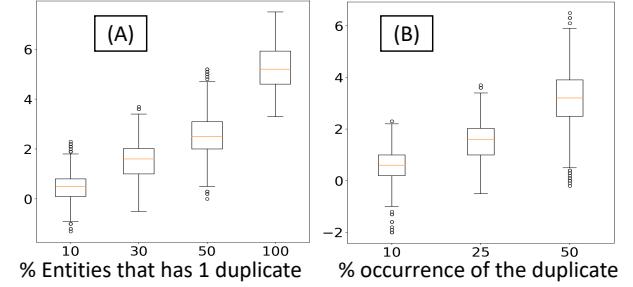


Figure 16: Hyperplane scenario results on HiCapRF with OHE. Only test set is diluted with duplicates. (A) Vary $|ED|/|E|$, while fixing $(|r|_t, occ(D_k), |D_k|) = (3000, 25, 1)$. (B) Vary $occ(D_k)$, while fixing $(|r|_t, |ED|/|E|, |D_k|) = (3000, 30, 1)$.

Y-axis: Delta drop in % test accuracy due to duplication with OHE

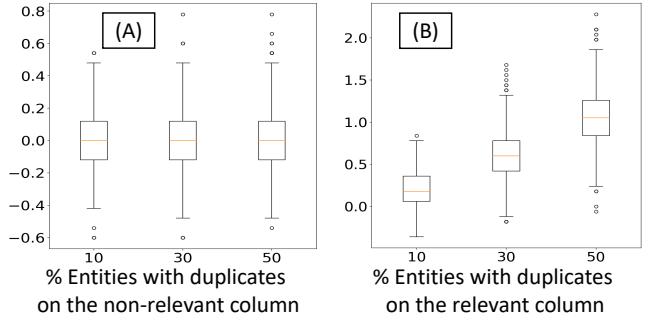


Figure 17: Hyperplane results on HiCapRF. We set $|\mathcal{A}| = 4$ and vary $|ED|/|E|$, while fixing $(|r|_t, occ(D_k), |D_k|) = (5000, 25, 1)$. Duplicates introduced on the column with (A) non-positive Relevancy (noisy column) (B) high Relevancy (predictive column).

11.2.2 Varying properties of duplicates being mapped to “Others” and column Relevancy. We now repeat the simulation scenario presented in Section 7.3.3 and Section 7.3.4 but with Hyperplane setting, i.e. the true distribution is given by a hyperplane that separates the classes. Figure 16 presents the results when only the test set is diluted with duplicates. We again note that the presence of duplicates in the test set impacts HiCapRF significantly. Figure 17 presents the Hyperplane simulation results when we have the presence of both high and low relevancy columns in the dataset (setup same as Section 7.3.4). We again find that the duplication on a noisy column has a marginal impact on HiCapRF, while duplicates on relevant column affect it significantly.

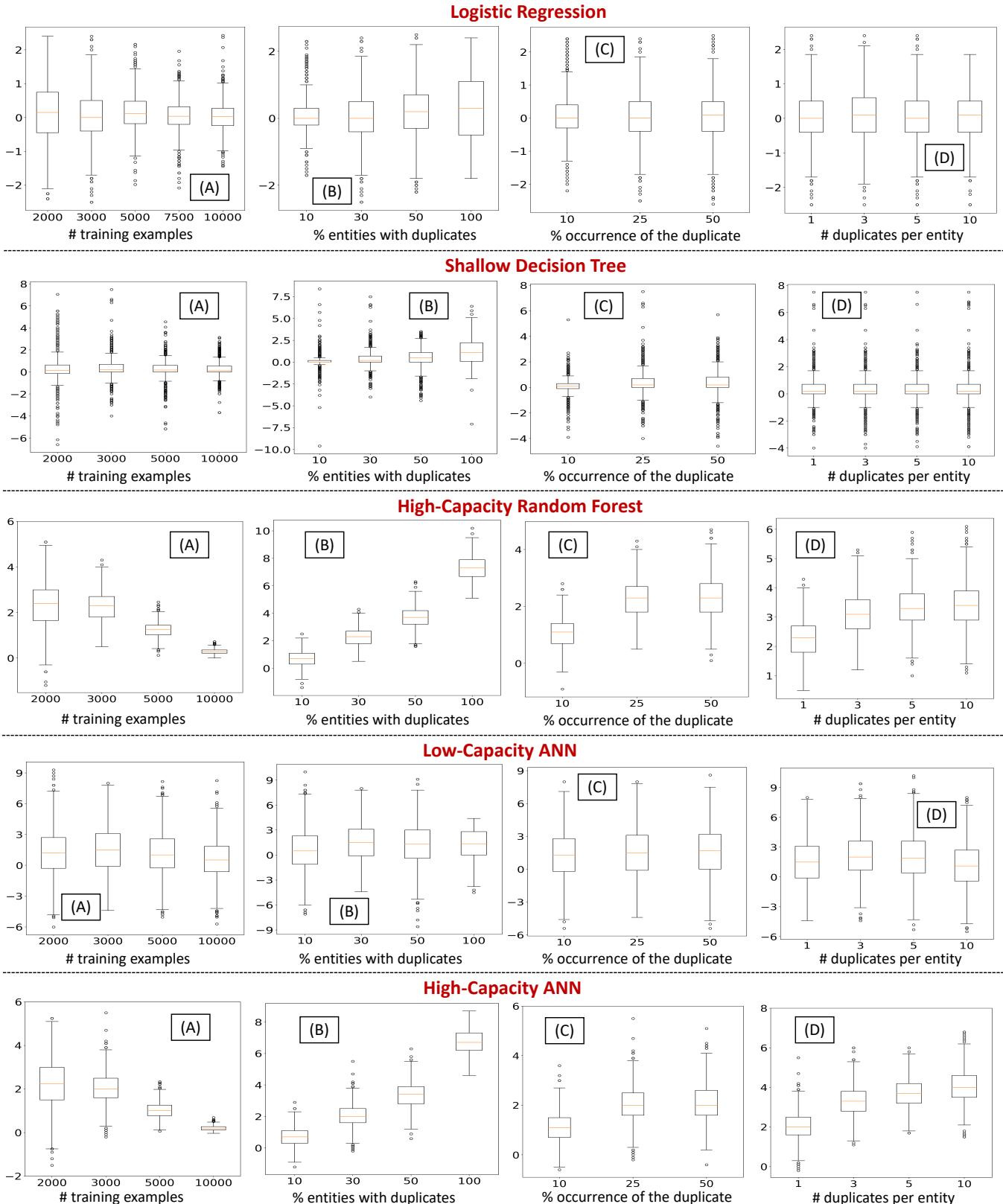


Figure 11: AllA simulation scenario results for LR, ShallowDT, HiCapRF, LoCapANN, and HiCapANN with OHE. $|X| = 3$. (A) Vary $|r|_t$ (# training examples), while fixing $(|ED|/|E|, occ(D_k), |D_k|) = (30, 25, 1)$. (B) Vary $|ED|/|E|$, while fixing $(|r|_t, occ(D_k), |D_k|) = (3000, 25, 1)$. (C) Vary $occ(D_k)$, while fixing $(|r|_t, |ED|/|E|, |D_k|) = (3000, 30, 1)$. (D) Vary $|D_k|$, while fixing $(|ED|/|E|, |r|_t, occ(D_k)) = (30, 3000, 25)$, for all $k \in [1, |ED|]$.

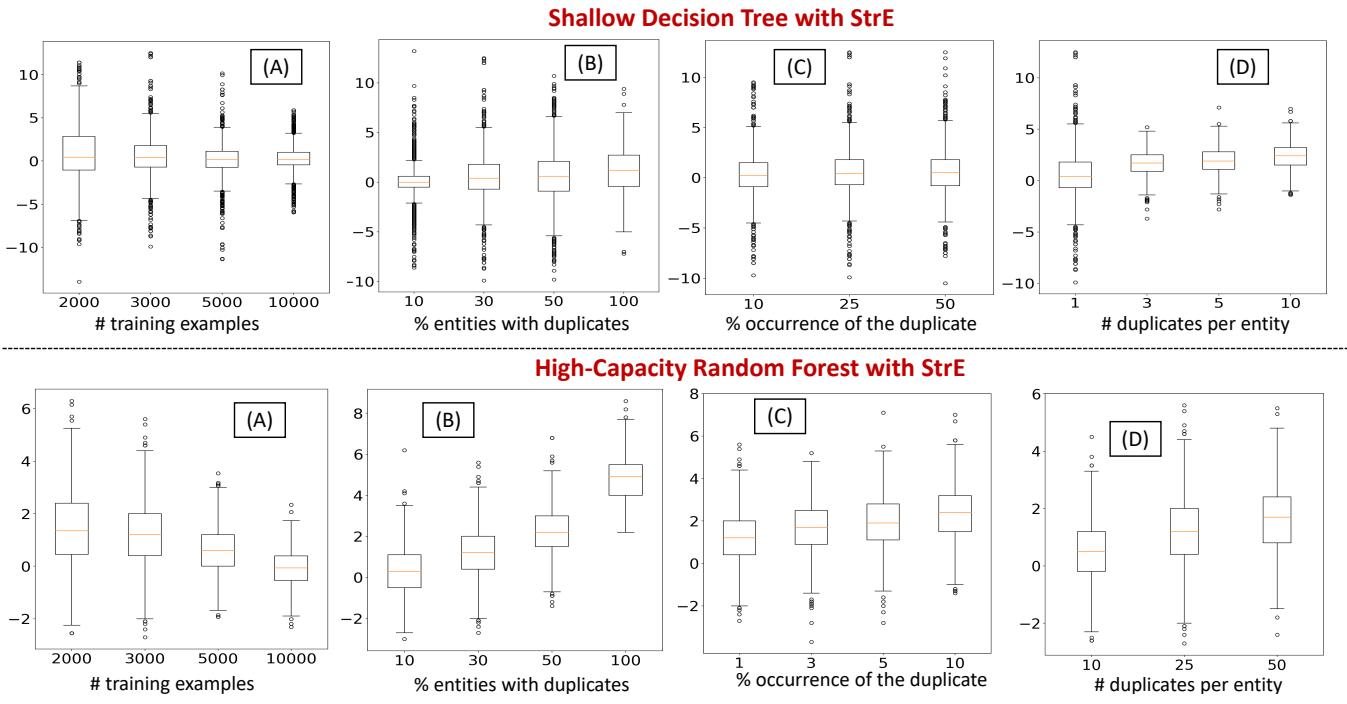


Figure 12: A11A simulation scenario results for ShallowDT and HiCapRF with StrE. (A) Vary $|r|_t$ (# training examples), while fixing $(|ED|/|E|, occ(D_k), |D_k|) = (30, 25, 1)$. (B) Vary $|ED|/|E|$, while fixing $(|r|_t, occ(D_k), |D_k|) = (3000, 25, 1)$. (C) Vary $occ(D_k)$, while fixing $(|r|_t, |ED|/|E|, |D_k|) = (3000, 30, 1)$. (D) Vary $|D_k|$, while fixing $(|ED|/|E|, |r|_t, occ(D_k)) = (30, 3000, 25)$, for all $k \in [1, |ED|]$.

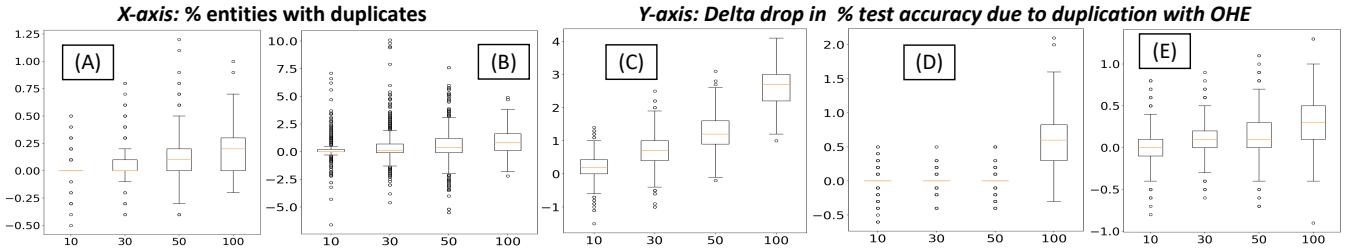


Figure 13: Hyperplane setting results on (A) LR (B) ShallowDT (C) HiCapRF (D) LoCapANN (E) HiCapANN (same setup as Figure 3(B)).

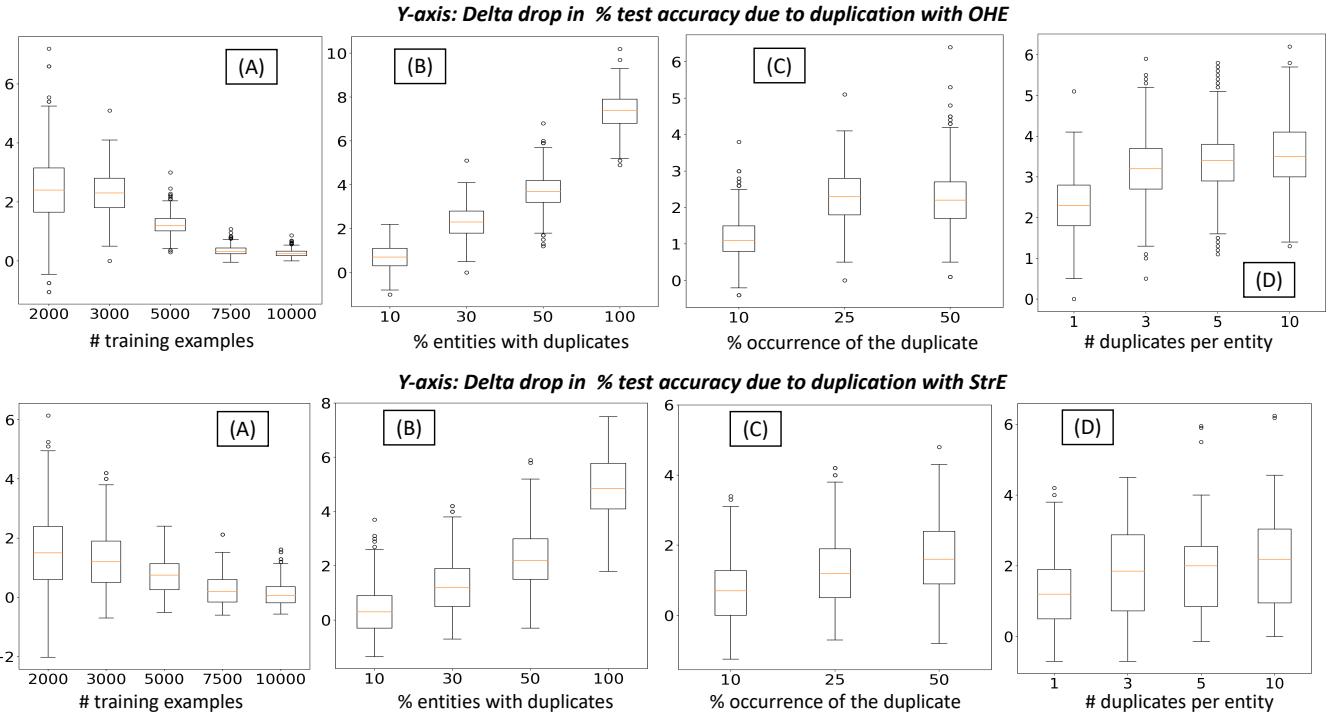


Figure 14: All A simulation scenario results for Random Forest with OHE and StrE where hyper-parameters are tuned with grid search. $|X| = 3$. (A) Vary $|r|_t$ (# training examples), while fixing $(|ED|/|E|, \text{occ}(D_k), |D_k|) = (30, 25, 1)$. (B) Vary $|ED|/|E|$, while fixing $(|r|_t, \text{occ}(D_k), |D_k|) = (3000, 25, 1)$. (C) Vary $\text{occ}(D_k)$, while fixing $(|r|_t, |ED|/|E|, |D_k|) = (3000, 30, 1)$. (D) Vary $|D_k|$, while fixing $(|ED|/|E|, |r|_t, \text{occ}(D_k)) = (30, 3000, 25)$, for all $k \in [1, |ED|]$.

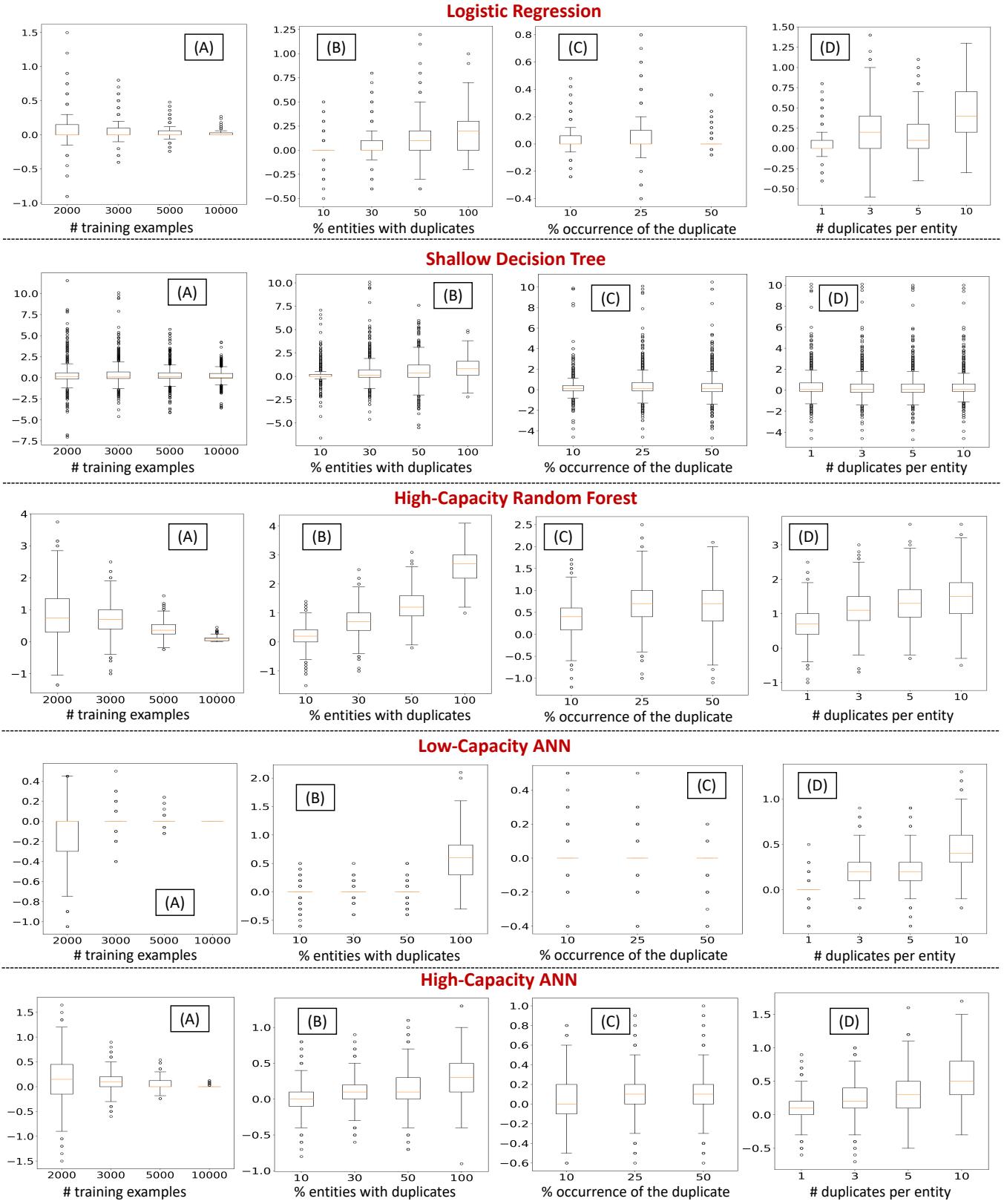


Figure 15: Hyperplane simulation scenario results for LR, ShallowDT, HiCapRF, LoCapANN, and HiCapANN with OHE. $|X| = 3$. (A) Vary $|r|_t$ (# training examples), while fixing $(|ED|/|E|, occ(D_k), |D_k|) = (30, 25, 1)$. (B) Vary $|ED|/|E|$, while fixing $(|r|_t, occ(D_k), |D_k|) = (3000, 25, 1)$. (C) Vary $occ(D_k)$, while fixing $(|r|_t, |ED|/|E|, |D_k|) = (3000, 30, 1)$. (D) Vary $|D_k|$, while fixing $(|ED|/|E|, |r|_t, occ(D_k)) = (30, 3000, 25)$, for all $k \in [1, |ED|]$.