

APPENDICES

Appendix 1 - Summary of Well-known Data Quality Issues and the Data Quality Dimension Violations Associated

Table 1 - Data Quality Issues and Data Quality Dimension Violations (Adapted from Visengeriyeva and Abedjan, 2020[31])

#	Data Quality Issue	Description	Data Quality Dimensions
1	Missing data	Comprises missing tuples and missing values. Tuple completeness requires that all tuples are present in the table. Missing value issue consists of either null values or disguised values. Value completeness requires that all values are present in the table, while null values indicate that the value is unknown or non-existent.	Accuracy, Completeness
2	Incorrect data	Data that differ from the values of the real-world entity (e.g., wrong date of birth).	Accuracy
3	Misspellings	Syntactic deviation of the data value from its ground truth (e.g., "Smiht" instead of "Smith").	Accuracy
4	Ambiguous data	Data values which might be interpreted in several ways (e.g., abbreviations or cryptic values).	Accuracy, Consistency
5	Extraneous data	Presence of additional data in the attribute value (e.g., the address column contains a person's name in addition to the address).	Consistency, Uniqueness
6	Outdated temporal data	Values that are obsolete or outdated.	Timeliness
7	Misfielded values	Values that are placed inside the wrong attribute.	Accuracy, Consistency, Completeness
8	Incorrect references	Entities that contain wrong information concerning the reference relation (e.g., the employee is associated with a wrong department).	Accuracy
9	Duplicates	Tuples/values that represent the same real-world entity.	Uniqueness
10	Structural conflicts	Conflicting duplicates in different sources.	Consistency, Uniqueness
11	Different word orderings	Values that violate the expected word order (e.g., first name precedes last name).	Consistency, Uniqueness
12	Different aggregation levels	Entities produced by applying different aggregation methods (e.g., entries per quartal vs. entries per year).	Accuracy, Consistency
13	Temporal mismatch	Refers to erroneous data that arise due to non-enforcement of integrity constraints for temporal data.	Accuracy, Timeliness
14	Different units/representations	Occurrence of multiple representations for the same concept (e.g., Price in different currencies).	Consistency
15	Domain violation	Values that violate semantic rules defined on the specific attribute.	Accuracy
16	FD violation	Values that violate previously specified functional dependencies.	Accuracy, Consistency
17	Wrong data type	Values that violate the data type specification of the corresponding attribute, i.e., data type constraint violation.	Consistency
18	Referential integrity violation	Tuples that violate the referential integrity constraints defined on multiple relations (e.g., missing foreign key).	Accuracy, Consistency, Completeness
19	Uniqueness violation	Duplication of values under the uniqueness constraint.	Uniqueness
20	Use of synonyms	Occurrence of synonymous representations for the same concept inside the same column (e.g., "lecturer" and "professor").	Uniqueness
21	Use of special characters (space, no space, dash, parentheses)	Refers to different representations of compound data, such as Social Security Number or phone number.	Consistency
22	Different encoding formats	Inconsistent usage of encodings for values within a dataset (e.g., ASCII or EBCDIC).	Consistency

Appendix 2 - Explanation for selection of 50 database tables from the Prague Relational Learning Repository [21]

These are the steps taken to obtain 50 database tables from the Prague Relational Learning Repository.

1. The paper that described the Repository (**Motl & Schulte (2024) - [21]**) had this Table 1:

Database	#Relations	#Instances	Size	Type	Domain	Task	In it we can see, in the first line, important information that allow us to choose the Databases. There are 49 Databases described here, some with numbers that point to research papers associated.
Accidents [18]	4	495760	210.0	Real	Gouverment	Class	
AdventureWorks	71	30669	234.6	Synth	Retail	Regr	
AustralianFootball	4	3036	38.0	Real	Sport	Class	
Biodegradability [3]	5	328	3.2	Real	Medicine	Regr	
Carcinogenesis [24]	6	329	26.3	Real	Medicine	Class	
CCS	6	1000	658.4	Real	Finance	Regr	
ClassicModels	8	273	0.5	Synth	Retail	Regr	
Countries	4	247	8.6	Real	Geography	Regr	
Credit	9	10084	443.6	Synth	Retail	Class	
CS	8	100	0.3	Synth	Finance	Class	
Dunur [11]	20	276	0.8	Real	Kinship	Class	
Elti [11]	14	1081	0.7	Real	Kinship	Class	
Employee	7	2838426	344.6	Synth	Retail	Regr	
Financial [2]	8	682	94.1	Real	Finance	Class	
FTP	2	29555	7.5	Synth	Retail	Class	
Genes [6]	3	862	1.9	Real	Medicine	Class	
Hepatitis	7	500	2.2	Real	Medicine	Class	
Hockey [23]	23	7759	15.5	Real	Sport	Class	
IMDb [23]	7	794625	614.6	Real	Entertainment	Class	
MovieLens [22]	7	6039	151.9	Real	Entertainment	Class	
Lahman	25	23111	84.0	Real	Sport	Regr	
LegalActs	5	564268	238.2	Real	Gouverment	Class	
Mesh [9]	32	223	1.1	Real	Industry	Regr	
Mondial [14]	33	454	3.3	Real	Geography	Class	
MooneyFamily [20]	72	92	3.3	Synth	Kinship	Class	
Mutagenesis [8]	3	188	0.9	Real	Medicine	Class	
Nations	3	14	2.1	Real	Geography	Class	
NBA [23]	4	30	0.3	Real	Sport	Class	
NCAA	10	268	40.6	Real	Sport	Class	
Northwind	29	830	1.1	Synth	Retail	Regr	
Pima	14	768	0.8	Real	Medicine	Class	
PremiereLeague [19]	4	363	11.3	Real	Sport	Class	
PTE [25]	41	299	7.3	Real	Medicine	Class	
Pubs	11	18	0.4	Synth	Retail	Regr	
Sakila	16	15991	6.6	Synth	Retail	Regr	
SalesDB	4	6148886	539.3	Synth	Retail	Regr	
SameGen	7	1081	0.3	Real	Kinship	Class	
Stats	8	38357	621.4	Real	Education	Regr	
StudentLoan [17]	13	1000	0.9	Real	Education	Class	
PTC [10]	4	343	7.8	Real	Medicine	Class	
Thrombosis [7]	3	806	1.9	Real	Medicine	Class	
TPCC [12]	9	28433	174.1	Synth	Retail	Class	
TPCDS [16]	24	99550	4587.5	Synth	Retail	Class	
TPCH [4]	8	148255	1925.1	Synth	Retail	Regr	
Trains [15]	2	20	0.1	Synth	Logistic	Class	
University	5	38	0.3	Synth	Education	Class	
UW-CSE [21]	4	278	0.2	Real	Education	Class	
VOC	8	8215	2.7	Real	Logistic	Class	
World	3	239	0.8	Real	Geography	Class	

Table 1: List of databases in the repository

2. The next step was to take a look at the tables themselves at the Prague Relational Learning Repository. We used MySQL Workbench to connect with the following parameters described at [21]:

Hostname: db.relationaldata.org

Port: 3306

Username: guest

Password: relational
3. After we understood the data we decided to extract the metadata, as we had done previously from UCI, including new information about Primary and Foreign Keys which are available as the tables are relational.
4. The next step was to download the information available then. There were over 2000 tables at that time. We just wanted 50.
5. We brought the following information, similar to what was brought earlier from all the 622 datasets of the UCI Machine Learning Repository:

index	name	database	instances	attributes	Column	Description	primary_key	foreign_keys
Accidents.nesreca	nesreca	Accidents	508993	22	<code>id_nesreca, klas_nesreca: id_nesreca (char(6)), klas_nesre_id_nesreca</code>			<code>upravna_enota -> upravna_enota.id_upravna_enota</code>
Accidents.oseba	oseba	Accidents	954036	17	<code>id_nesreca, povrocitelj_id_nesreca (char(6)), povrocitelj_all_udelenece (id_nesreca -> nesreca.id_nesreca, upravna_enota -> upravna_enota)</code>			
Accidents.upravna_enota	upravna_enota	Accidents	64	4	<code>id_upravna_enota, ime_id_upravna_enota (char(4)), im_id_upravna_enota</code>			
AdventureWorks2014.AWBuildVersion	AWBuildVersion	AdventureWorks2014	1	4	<code>SystemInformationID, SystemInformationID (tinyint) or SystemInformationID</code>			
AdventureWorks2014.Address	Address	AdventureWorks2014	19614	9	<code>AddressID, AddressLine1AddressID (int), AddressLine1 (AddressID, Name, rAddressTypeID (int), Name (var AddressTypeID))</code>		<code>StateProvinceID -> StateProvince.StateProvinceID</code>	
AdventureWorks2014.AddressType	AddressType	AdventureWorks2014	6	4	<code>AddressTypeID, Name, rAddressTypeID (int), Name (var AddressTypeID)</code>			
AdventureWorks2014.BillOfMaterials	BillOfMaterials	AdventureWorks2014	2679	9	<code>BillOfMaterialsID, ProductBillOfMaterialsID (int), ProductID, ComponentID -> Product.ProductID, ProductAssembly</code>			

- a. The index, with the content of the database followed by a dot and the table
 - b. The name of the table
 - c. The name of the database
 - d. The number of Instances/records in this table
 - e. The number of attributes/columns in this table
 - f. The columns in this table
 - g. The description, where the columns contain the format associated
 - h. The primary_key(s) in this table
 - i. The foreign_key(s) in this table
6. From Table 1 of paper [21] we checked if the databases existed in the current repository. Unfortunately, some databases did not exist anymore or had different names at the time.
7. These were the 3 databases not found then: PTC, Thrombosis and VOC.
8. Analyzing the content of all different Domains it was discovered that the Kinship Domain with 4 databases, was not suitable for our research, because all their tables contain only few columns with information about names of relatives, such as brother, father, etc., so nothing usable for our analysis:

Dunur.brother		Dunur.husband		Elti.target		
name1	name2	name1	name2	name1	name2	is_elti
Alfonso	Sophia	Andrew	Christine	ali1	ali2	0
Arthur	Victoria	Arthur	Margaret	ali1	alp	0
Colin	Charlotte	Charles	Jennifer	ali1	anil	0
Emilio	Lucia	Christopher	Penelope	ali1	ayse	0
James	Jennifer	Emilio	Gina	ali1	ayten	0
Marco	Angela					

9. So, the initial total of 49 databases became $49 - 3$ (not found) – 4 (Kinship) = 42 databases.

10. See below the list of the 49 databases, with the Current Name and the Name in Table 1:

Current Name	Name in Table 1	#Relations	#Instances	Size(MB)	Type	Domain	Task
Accidents	Accidents	4	495760	210	Real	Government	Class
AdventureWorks2014	AdventureWorks	71	30669	234.6	Synth	Retail	Regr
AustralianFootball	AustralianFootball	4	3036	38	Real	Sport	Class
Biodegradability	Biodegradability	5	328	3.2	Real	Medicine	Regr
Carcinogenesis	Carcinogenesis	6	329	26.3	Real	Medicine	Class
ccs	CCS	6	1000	658.4	Real	Finance	Regr
classicmodels	ClassicModels	8	273	0.5	Synth	Retail	Regr
Countries	Countries	4	247	8.6	Real	Geography	Regr
Credit	Credit	9	10084	443.6	Synth	Retail	Class
cs	CS	8	100	0.3	Synth	Finance	Class
Dunur	Dunur	20	276	0.8	Real	Kinship	Class
Elti	Elti	14	1081	0.7	Real	Kinship	Class
employee	Employee	7	2838426	344.6	Synth	Retail	Regr
financial	Financial	8	682	94.1	Real	Finance	Class
ftp	FTP	2	29555	7.5	Synth	Retail	Class
genes	Genes	3	862	1.9	Real	Medicine	Class
Hepatitis_std	Hepatitis	7	500	2.2	Real	Medicine	Class
Hockey	Hockey	23	7759	15.5	Real	Sport	Class
imdb_full	IMDb	7	794625	614.6	Real	Entertainment	Class
imdb_MovieLens	MovieLens	7	6039	151.9	Real	Entertainment	Class
lahman_2014	Lahman	25	23111	84	Real	Sport	Regr
legalActs	LegalActs	5	564268	238.2	Real	Government	Class
Mesh	Mesh	32	223	1.1	Real	Industry	Regr
Mondial	Mondial	33	454	3.3	Real	Geography	Class
MooneyFamily	MooneyFamily	72	92	3.3	Synth	Kinship	Class
mutagenesis	Mutagenesis	3	188	0.9	Real	Medicine	Class
nations	Nations	3	14	2.1	Real	Geography	Class
NBA	NBA	4	30	0.3	Real	Sport	Class
NCAA	NCAA	10	268	40.6	Real	Sport	Class
northwind	Northwind	29	830	1.1	Synth	Retail	Regr
Pima	Pima	14	768	0.8	Real	Medicine	Class
PremierLeague	PremiereLeague	4	363	11.3	Real	Sport	Class
Database not found	PTC	1	343	7.8	Real	Medicine	Class
PTE	PTE	41	299	7.3	Real	Medicine	Class
pubs	Pubs	11	18	0.4	Synth	Retail	Regr
sakila	Sakila	16	15991	6.6	Synth	Retail	Regr
SalesDB	SalesDB	4	6148886	539.3	Synth	Retail	Regr

Same Gen	SameGen	1	1081	0.3	Real	Kinship	Class
stats	Stats	8	38357	621.4	Real	Education	Regr
Student_loan	StudentLoan	13	1000	0.9	Real	Education	Class
Database not found	Thrombosis	1	806	1.9	Real	Medicine	Class
tpcc	TPCC	9	28433	174.1	Synth	Retail	Class
tpcds	TPCDS	24	99550	4587.5	Synth	Retail	Class
tpch	TPCH	8	148255	1925.1	Synth	Retail	Regr
trains	Trains	2	20	0.1	Synth	Logistic	Class
university	University	5	38	0.3	Synth	Education	Class
UW_std	UW-CSE	4	278	0.2	Real	Education	Class
Database not found	VOC	1	9215	2.7	Real	Logistic	Class
world	World	3	239	0.8	Real	Geography	Class

11. An analysis was made, and this was the distribution of the databases:

Category	Subcategory	Count
Domain	Retail	12
	Medicine	7
	Sport	6
	Geography	4
	Education	4
	Kinship	0
	Finance	3
	Government	2
	Entertainment	2
	Industry	1
	Logistic	1
	Total	42
Size	Small (< 1,000 instances)	16
	Medium (1,000 - 100,000 instances)	18
	Large (> 100,000 instances)	8
	Total	42
Type	Real	28
	Synthetic	14
	Total	42
Task	Classification	25
	Regression	17
	Total	42

12. The next step was filtering to 25 databases. It was found out that the Accidents' Database tables had Czech words in titles, not English words, which means that our code could not check the words in them. So, one more database was cancelled. It was from the Government Domain. A proportion was made regarding the Count of Domain databases:

Category	Subcategory	Count	Filtered Count
Domain	Retail	12	7
	Medicine	7	4
	Sport	6	3
	Geography	4	2
	Education	4	2
	Kinship	0	0
	Finance	3	2
	Government	2	1
	Entertainment	2	2
	Industry	1	1
	Logistic	1	1
	Total	42	25

13. This was the list of the 25 databases chosen, by alphabetical order of Domain

Current Name	Name in Table 1	#Relations	#Instances	Size (MB)	Type	Domain	Task
stats	Stats	8	38357	621.4	Real	Education	Regr
UW_std	UW-CSE	4	278	0.2	Real	Education	Class
imdb_full	IMDb	7	794625	614.6	Real	Entertainment	Class
imdb_MovieLens	MovieLens	7	6039	151.9	Real	Entertainment	Class
ccs	CCS	6	1000	658.4	Real	Finance	Regr
financial	Financial	8	682	94.1	Real	Finance	Class
Countries	Countries	4	247	8.6	Real	Geography	Regr
Mondial	Mondial	33	454	3.3	Real	Geography	Class
legalActs	LegalActs	5	564268	238.2	Real	Government	Class
Mesh	Mesh	32	223	1.1	Real	Industry	Regr
trains	Trains	2	20	0.1	Synth	Logistic	Class
Biodegradability	Biodegradability	5	328	3.2	Real	Medicine	Regr
Carcinogenesis	Carcinogenesis	6	329	26.3	Real	Medicine	Class
Hepatitis_std	Hepatitis	7	500	2.2	Real	Medicine	Class
PTE	PTE	41	299	7.3	Real	Medicine	Class
classicmodels	ClassicModels	8	273	0.5	Synth	Retail	Regr
Credit	Credit	9	10084	443.6	Synth	Retail	Class
ftp	FTP	2	29555	7.5	Synth	Retail	Class
northwind	Northwind	29	830	1.1	Synth	Retail	Regr
sakila	Sakila	16	15991	6.6	Synth	Retail	Regr
SalesDB	SalesDB	4	6148886	539.3	Synth	Retail	Regr
tpcds	TPCDS	24	99550	4587.5	Synth	Retail	Class
AustralianFootball	AustralianFootball	4	3036	38	Real	Sport	Class
Hockey	Hockey	23	7759	15.5	Real	Sport	Class
Lahman_2014	Lahman	25	23111	84	Real	Sport	Regr

14. This list of 25 databases contained 331 tables. The code used to download the information, 'BringingDatabaseInformation.ipynb' and the resulting file created 'all_selected_databases_info.xlsx' are available on our GitHub [29].
 15. The final definition was based in analysis of the content, to choose interesting tables from the maximum amount of the 25 databases.
 16. Some databases did not contain interesting information for attribute analysis, such as Carcinogenesis and Mesh, and were discarded.
 17. Another decision was made regarding the number of Foreign Keys. In a future step of our research, we will analyse the relationship between different tables, and the fact that there are tables with many FKs will enable that, and this was an important decision taken.
 18. This is the final definition of the fifty tables:

This final definition is in the file Fifty_Database_Tables_from_Prague_Repository.xlsx [29]

This is the analysis of the fifty final tables:

Database Distribution Analysis		Total number of databases: 23
		Area Distribution Analysis
Current Name	Number	
sakila	10	• Retail: 21 tables
Mondial	5	• Sport: 7 tables
northwind	5	• Entertainment: 6 tables
lahman_2014	4	• Geography: 6 tables
imdb_full	3	• Medicine: 4 tables
imdb_MovieLens	3	• Education: 2 tables
AustralianFootball	2	• Finance: 2 tables
PTE	2	• Government: 1 table
tpcds	2	• Logistic: 1 table
Biodegradability	1	
ccs	1	
classicmodels	1	
Countries	1	
Credit	1	
financial	1	
ftp	1	
Hepatitis_std	1	
Hockey	1	
legalActs	1	
SalesDB	1	
stats	1	
trains	1	
UW_std	1	
Carcinogenesis	0	
Mesh	0	
Total	50	Total number of areas: 9

Table Size Distribution Analysis	
• Small (0-1,000 instances): 23 tables	
• Medium (1,001-100,000 instances): 20 tables	
• Large (100,001+ instances): 7 tables	

Notable large tables:	
1.	Credit.charge: 1,600,000 instances
2.	financial.trans: 1,056,320 instances
3.	SalesDB.Sales: 6,715,221 instances
4.	tpcds.store_sales: 2,880,404 instances

Primary Key Analysis

Total tables: 50 Tables

Tables with defined Primary Keys: 47 Tables (39 use single-column primary keys and 8 use composite keys).

3 Tables without defined Primary Keys:

1. Hockey.Master
2. northwind.Invoices
3. northwind.Order Details Extended

Composite Key Analysis

Total tables with composite keys: 8

Tables with composite keys:

1. imdb_full.movies2actors (movieid, actorid)
 2. imdb_MovieLens.u2base (userid, movieid)
 3. lahman_2014.halloffame (hofID, yearID)
 4. lahman_2014.salaries (yearID, teamID, lgID, playerID)
 5. lahman_2014.teams (yearID, lgID, teamID)
 6. Mondial.city (Name, Province)
 7. Mondial.province (Name, Country)
 8. tpcds.store_sales (ss_item_sk, ss_ticket_number)
- The lahman_2014 database uses composite keys most frequently, with 3 out of its 4 tables having composite keys.
 - The Mondial database uses composite keys for geographic entities, likely to handle cases where names might be repeated across different regions.
 - The imdb databases use composite keys for relationship tables, which is a common practice in many-to-many relationships.

Foreign Keys analysis

Total tables: 50 Tables

Tables with defined Foreign Keys: 26 Tables

index	#FKs	
tpcds.store_sales	9	
Credit.charge	4	
Mondial.city	3	
northwind.Orders	3	
sakila.rental	3	
SalesDB.Sales	3	
imdb_full.movies2actors	2	
imdb_MovieLens.u2base	2	
Mondial.politics	2	
northwind.Products	2	
sakila.customer	2	
sakila.film	2	
sakila.inventory	2	
sakila.staff	2	
sakila.store	2	
classicmodels.customers	1	
financial.trans	1	
Hockey.Master	1	
lahman_2014.players	1	
lahman_2014.salaries	1	
Mondial.population	1	
Mondial.province	1	
PTE.pte_atm	1	
sakila.address	1	
sakila.city	1	
trains.cars	1	

Observe that the databases with more tables chosen contain FKs, such as sakila, Mondial, and northwind, allowing future FK analysis:

Current Name	Number
sakila	10
Mondial	5
northwind	5
lahman 2014	4
imdb full	3
imdb MovieLens	3
AustralianFootball	2
PTE	2
tpcds	2

Appendix 3 – Extraction of Abbreviation–Full Term Pairs from the NameGuess Dataset

To enhance the semantic type detection process, the abbreviations dictionary was expanded by leveraging evaluation data from the *NameGuess* project [37]. The *NameGuess* system, introduced at EMNLP 2023, focuses on reducing data quality risks from cryptic column headers through bidirectional transformation between abbreviated and expanded column names.

1. Repository and Dataset Structure

The *NameGuess* GitHub repository (<https://github.com/amazon-science/nameguess>) provides a collection of JSON files representing evaluation datasets. These files are located in a structured subdirectory (e.g., nameguess/data/) and each JSON object contains:

- `table_id`: Identifier of the source table,
- `technical_name`: The abbreviated column name,
- `gt_label`: The expanded, ground truth label,
- `difficulty`: Optional metadata field.

2. Parsing and Extraction Process

A Python notebook file ('**Obtaining Items from Name Guess.ipynb**' [29]) was developed to iterate through the relevant JSON files (eval_chicago.json, eval_la.json and eval_sf.json). The extraction followed these steps:

- a. **Selection**: JSON files containing the substring "eval" in their names were selected for processing.
- b. **Parsing**: The script used the json module to load each file and iterate through entries.
- c. **Matching**: For each entry, `technical_name` was compared to `gt_label`. If the two differed and a valid abbreviation–full term pattern was found, the pair was extracted.
- d. **Aggregation**: Each extracted pair was recorded along with the file name and table ID. The collected data were exported into a CSV file for inspection.

3. Dictionary Expansion and Integration

The resulting list of abbreviation–full form pairs was reviewed, deduplicated, and integrated with the existing dictionary. This increased the size of the abbreviations dictionary from 300 to **839 entries**, significantly enhancing coverage of commonly used column header abbreviations.

These new abbreviations also contributed to expanding the *Formats Dictionary*, particularly for semantic types commonly represented through abbreviations (e.g., 'hrs' → 'hours', 'amt' → 'amount', 'dob' → 'date_of_birth'). This update improved the accuracy of rule-based semantic type detection, particularly in real-world data sources with non-standardised column labels.

See some interesting abbreviations:

10th: tenth
abbrev: abbreviation
abdnnt: abandonment
abndnd: abandoned
abs: absolute
admin: administrator
advsry: advisory
aero: aerospace
afdble: affordable
affiliatn: affiliation
affirmatn: affirmation
affltd: affiliated