

Automated Detection of Records in Biological Sequence Databases that are Inconsistent with the Literature

Mohamed Reda Bouadjenek, Karin Verspoor, Justin Zobel

Department of Computing and Information Systems, University of Melbourne, Parkville, 3053, Australia

Abstract

We investigate and analyse the data quality of nucleotide sequence databases with the objective of automatic detection of data anomalies and suspicious records. Specifically, we demonstrate that the published literature associated with each data record can be used to automatically evaluate its quality, by cross-checking the consistency of the key content of the database record with the referenced publications. Focusing on GenBank, we describe a set of quality indicators based on the relevance paradigm of information retrieval (IR). Then, we use these quality indicators to train an anomaly detection algorithm to classify records as “*confident*” or “*suspicious*”.

Our experiments on the PubMed Central collection show assessing the coherence between the literature and database records, through our algorithms, is an effective mechanism for assisting curators to perform data cleansing. Although fewer than 0.25% of the records in our data set are known to be faulty, we would expect that there are many more in GenBank that have not yet been identified. By automated comparison with literature they can be identified with a precision of up to 10% and a recall of up to 30%, while strongly outperforming several baselines. While these results leave substantial room for improvement, they reflect both the very imbalanced nature of the data, and the limited explicitly labelled data that is available. Overall, the obtained results show promise for the development of a new kind of approach to detecting low-quality and suspicious sequence records based on literature analysis and consistency. From a practical point of view, this will greatly help curators in identifying inconsistent records in large-scale sequence databases by highlighting records that are likely to be inconsistent with the literature.

Keywords: Data Analysis, Data Quality, Bioinformatics Databases, Anomaly Detection.

1. Introduction

Bioinformatics sequence databases such as GenBank or UniProt contain large numbers of nucleic acid sequences and protein sequences. In 2016, GenBank alone contained over 224 billion nucleotide bases in more than 198 million sequences — a number that is growing at an exponential rate, doubling every 18 months.¹

Email addresses: reda.bouadjenek@unimelb.edu.au (Mohamed Reda Bouadjenek), karin.verspoor@unimelb.edu.au (Karin Verspoor), jzobel@unimelb.edu.au (Justin Zobel)

¹<http://www.ncbi.nlm.nih.gov/GenBank/statistics/>

1
2
3
4 In commercial organizations, the primary reason for creating and maintaining such databases is their impor-
5 tance in the process of drug discovery, while in research they are used to understand the biological basis of
6 disease. Thus, a high level of data quality is crucial.
7

8
9 However, since these databases are fed by direct submissions from individual laboratories and by bulk
10 submissions from large-scale sequencing centers, they suffer from a range of data quality issues ([29]) including
11 errors, redundancies, ambiguities, incompleteness, and as we will show, discrepancies such as inconsistency
12 with the literature. Most of these records are linked to research articles in which the sequence was reported,
13 but the need to manually create the records on such a large scale means that errors creep in and, given the
14 volume, human curation alone is not sufficient for detection of these errors.
15

16
17 In this work, we seek to investigate and analyse the data quality of sequence databases from the perspec-
18 tive of a curator, who must detect anomalous and suspicious records. In contrast to previous research, which
19 has concerned detection of duplicate records ([12, 30]) and erroneous annotations ([7, 26, 43]), we emphasize
20 detection of low-quality records that we define as being inconsistent with the published literature. Specifi-
21 cally, we propose that the literature that is linked to records in their “reference” fields be automatically used
22 as background knowledge to check their quality. We explore a combination of information retrieval (IR) and
23 machine learning techniques to identify records that are anomalous and thus merit analysis by a curator.
24

25
26 To provide insight into the data quality of the nucleotide records cited by articles available in PubMed
27 Central² (PMC) from a literature consistency point of view, we analyzed these records as illustrated in
28 Figure 1. This figure shows the term overlap similarity³ between the record definition and different sections
29 of its associated article(s) (representing the title, abstract, body, and the full text). There are three notable
30 trends here: first, term overlap increases from title to body and full text since the size grows accordingly;
31 second, there is a high term overlap of roughly 80% between the record description field and the literature
32 body section; and third, for a small number of records, in which the overlap similarity is below 0.2, there is
33 low overlap or no overlap at all between the description field and the full text of their associated articles, thus
34 statistically suggesting a data quality problem. As an example, the record with accession number KM403369⁴
35 doesn’t share any terms with the article PMC4465667⁵ that is supposed to report on that record. Compared
36 to the median value, which is roughly 80% similarity between a record description field and the body section
37 of the article (see Figure 1), this association can be considered an outlier from a statistical perspective,
38 and can be argued to be weak. While this observation is purely statistical, it may be an indicator of a low
39 confidence in that record. Although this record is not necessarily faulty, its characteristics in relation to the
40 overall statistical distribution clearly suggest that it should be flagged as “suspicious”, and should be sent
41 to a curator for further investigation.
42
43
44
45
46
47
48
49
50
51
52
53
54
55

56 ²<http://www.ncbi.nlm.nih.gov/pmc/>

57 ³We use the overlap similarity to emphasize the number of terms of a record definition that are in its associated article.

58 Here, $Overlap(X_1, X_2) = |X_1 \cap X_2| / \min(|X_1|, |X_2|)$.

59 ⁴<http://www.ncbi.nlm.nih.gov/nucleotide/KM403369>

60 ⁵<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4465667/>

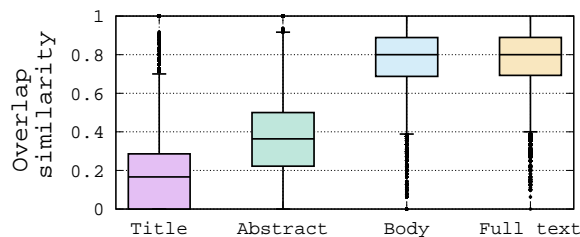


Figure 1: Overlap similarity between a record definition field and different sections of its associated document.

Usually, a suspicious record is reported manually, by a curator whose the job consists mainly to check the database records, the record's original submitter, or a third person who may use the database and notice the inconsistency of that record. To illustrate the difficulty of the task of identifying failing records, we analysed the distribution of record ages, for records which have been removed. This analysis showed that removed records have an average age of about 1 month at their removal time. This leads us to make two hypotheses: either (i) it takes about one month for a problematic record to be detected, or (ii) curators focus only on new records, while neglecting older ones. Either way, it is clear that there is a time window of only 1 month during which curators act. Hence, if a suspicious record is not identified in this time frame, it has a low probability of being spotted. These observations show the difficulty of the curator's job, and the need for the development of automatic methods to assist them.

With the aim of assisting curators, and while focusing on GenBank, we present in this paper a method for detection of suspicious records based on their associated articles and also on the collection of articles as a whole. To the best of our knowledge, this work is the first to use the literature for data quality assessment of bioinformatics sequence databases. The contributions of this paper are as follows:

- We demonstrate that the research literature can be automatically used for assessing the quality of a record.
- We propose a list of quality indicators that correlate with the quality of a record. The quality indicators are then used to train a learning anomaly detection algorithm.
- Our experiments on the full PubMed Central collection show that, although less than 0.25% of the records in our data set are faulty, by automated comparison with literature they can be identified with a precision of up to 10% and a recall of up to 30%, while greatly outperforming the best baseline.

2. Related Work

There is a substantial body of research related to data quality in bioinformatics databases. Previous research has focused mainly on duplicate record detection and erroneous annotations, as reviewed below.

2.1. Duplicate records

Koh et al. [30] use association rule mining to check for duplicate records with per-field exact, edit distance, or BLAST sequence [1] alignment matching. Drawbacks of this method, and its poor performance, have been shown by Chen et al. [12]. Similarly, Apiletti et al. [3] proposed extraction of association rules among attribute values to find causality relationships among them. By analysing the support and confidence of each rule, the method can show the presence of erroneous data. Other approaches also use approximate string matching to compute metadata similarity [46, 52, 10]. However, as they focus only on metadata, the underlying interpretation is that duplicates are assumed to have high metadata similarity, or that their sequences are identical.

Other approaches consider duplicates at the sequence level; they examine sequence similarity and use a similarity threshold to identify duplicates. For example, Holm and Sander [23] identified pairs of records with over 90% mutual sequence identity. Heuristics have been used in some of these methods to skip unnecessary pairwise comparisons, thus improving the efficiency. Li and Godzik [35] proposed CD-HIT, a fast sequence clustering method that uses heuristics to estimate the anticipated sequence identity and will skip the sequence alignment if the pair is expected to have low identity. Recently, Zorita et al. [60] proposed Star Code to detect duplicate sequences, which uses the edit distance as a threshold and will skip pairs exceeding the threshold. Such methods are valuable for this task, but do not address the problem of consistency or anomaly.

2.2. Erroneous annotations

Sequence databases exist as a resource for biomedicine, but the utility of the sequence of an organism depends on the quality of its annotations [10]. The annotations indicate the locations of genes and the coding regions in a sequence, and indicate what those genes do. That is, annotations serve as a reading guide to a sequence, which makes the scientific community highly reliant on this information. Although the research and development of algorithms for identifying coding sequences (CDSs) is still an active area in bioinformatics research, genome annotation has evolved greatly during past few years [55, 17, 15, 51]. However, the functional annotation of CDSs is particularly difficult to automate [57]. Current state-of-the-art functional annotation methods integrate multiple types of evidences [11, 4, 37], but unfortunately the quality of functional annotations remains generally poor [2, 56, 49] and is highly dependent on resource-intensive manual curation [40, 42].

Previous research work on function annotation identified potential problems with large-scale annotation efforts [7, 5, 41], and misannotation is a growing concern among the general research community, as mis-annotations can have a several impacts in diverse biological areas [44, 19, 25]. Even in very small bacterial genomes, many misannotations may arise [31]. As for high-throughput functional annotations, errors may occur due to a variety of factors [34, 39], but the most common errors are over-annotations, in which a gene is given a specific but incorrect function [49, 31, 36]. Once made, functional annotation errors can be difficult to correct in large scale sequence databases and as functional annotations are often inferred

1
2
3
4 from sequence similarity to other annotated sequences, errors may “propagate” to newly sequenced genomes
5 through “(mis)annotation transfer” [18, 24, 43].
6

7 Overall, existing data quality analysis methods for sequence databases focus only on the internal char-
8 acteristics of records. Our work demonstrate that the literature associated to records is a valuable external
9 source of information for assessing the quality of sequence database records.
10
11

12 13 14 **3. Background and Problem Definition**

15
16 In this section, we first provide an overview of GenBank, the most commonly used sequence database,
17 and we describe the structure of a sequence record in GenBank. Next, we discuss several kinds of data issues
18 in bioinformatics sequence databases, and finally, we define in detail the problem we study.
19
20

21 22 *3.1. GenBank overview*

23
24 In this work, we mainly focus on GenBank as it is the most important and most influential sequence
25 archive repository for research in almost all biological fields, whose data are accessed and cited by millions of
26 researchers around the world. GenBank is produced and maintained by the National Center for Biotechnology
27 Information (NCBI), and is considered to be an archive rather than a database, because multiple versions of
28 a given record may be maintained for historical purposes. The sequence submission to GenBank can occur
29 through: (i) direct submissions from scientists using BankIt⁶, which is a Web-based form, or the stand-alone
30 submission program, Sequin⁷, or (ii) bulk submissions most often done by large-scale sequencing centers,
31 which include Expressed Sequence Tag (EST), Sequence-tagged site (STS), Genome Survey Sequence (GSS),
32 and High-Throughput Genome Sequence (HTGS). Upon receipt of a sequence submission, an accession
33 number is assigned to the sequence, and then, it is released to the public database, where the entry is
34 retrievable using Entrez⁸.
35
36
37
38
39
40

41 Due to the fact that records can be submitted by multiple research actors without any particular data
42 quality control, errors may occur. Errors can seriously hamper the efficacy of analysis, data mining, and
43 machine learning algorithms. Hence, a faulty record is usually reported manually, by a database curator,
44 the record’s original submitter, or a third person who may use the database and notice the inconsistency of
45 that record. However, updates and revisions of a GenBank sequence can also be made by the submitters at
46 any time.
47
48
49

50 In addition to GenBank which can be considered as an unreviewed repository and thus may contain low
51 quality sequences, NCBI also maintains other curated sequence databases such as the Reference Sequence
52 (RefSeq). RefSeq provides a comprehensive, integrated, non-redundant, well-annotated set of sequences,
53
54

55
56
57 ⁶<https://www.ncbi.nlm.nih.gov/WebSub/?tool=GenBank>

58 ⁷<https://www.ncbi.nlm.nih.gov/Sequin/>

59 ⁸The Entrez Global Query Cross-Database Search System is a federated search engine, or web portal that allows users to
60 search many discrete health sciences databases at the National Center for Biotechnology Information (NCBI) website.
61
62
63
64
65

1
2
3
4 including genomic DNA, transcripts, and proteins. RefSeq genomes are copies of selected assembled genomes
5 available in GenBank, which have been generated by several processes including manual curation, which is
6 known to be a tedious and painful task.
7
8

9 10 *3.2. GenBank sequence record structure*

11 The format of a sequence record can be regarded as having three parts: the header, which contains the
12 information that applies to the whole record; the features, which are the annotations on the sequence; and
13 the sequence itself. The header section is composed of several fields:
14
15

- 16
17 • *LOCUS* field: contains a number of different data elements, including locus name, sequence length,
18 molecule type, and modification date. The locus name is designed to help group entries with similar
19 sequences: the first three characters usually designate the organism; the fourth and fifth characters
20 can be used to show other group designations, such as gene product; for segmented entries the last
21 character is one of a series of sequential integers;
22
- 23 • *DEFINITION* field: a brief description of sequence or sequence's function;
24
- 25 • *ACCESSION* field: a unique identifier for the record;
26
- 27 • *SOURCE* field: gives information about the sequence's organism;
28
- 29 • *REFERENCE* field: lists a set of publications by the authors of the sequence that discuss the data
30 reported in the record.
31
32

33 It is clear that the header part represents in itself a rich source of information.
34

35
36 Based on the fact that articles discuss the data reported in the records, and that there is high term
37 overlap between the record definition and its associated articles as reported in Figure 1, we will primarily
38 focus on the record definitions to assess the quality of the records from a literature consistency point of view.
39
40

41 42 *3.3. Classification of Biological Data Quality Issues*

43
44 Given a sequence record with its multiple data elements, the complex sequence submission process and
45 the data integration pipeline defined to exchange data with other sequence databases, data quality issues
46 may have physical or conceptual sources. Hence, Koh et al. [28, 29] proposed to classify biological data issues
47 according to their presence in data items at mainly four levels of detail — individual attributes, individual
48 records, individual databases, and multiple databases, as shown in Figure 2. Below, we briefly discuss these
49 data quality issues.
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

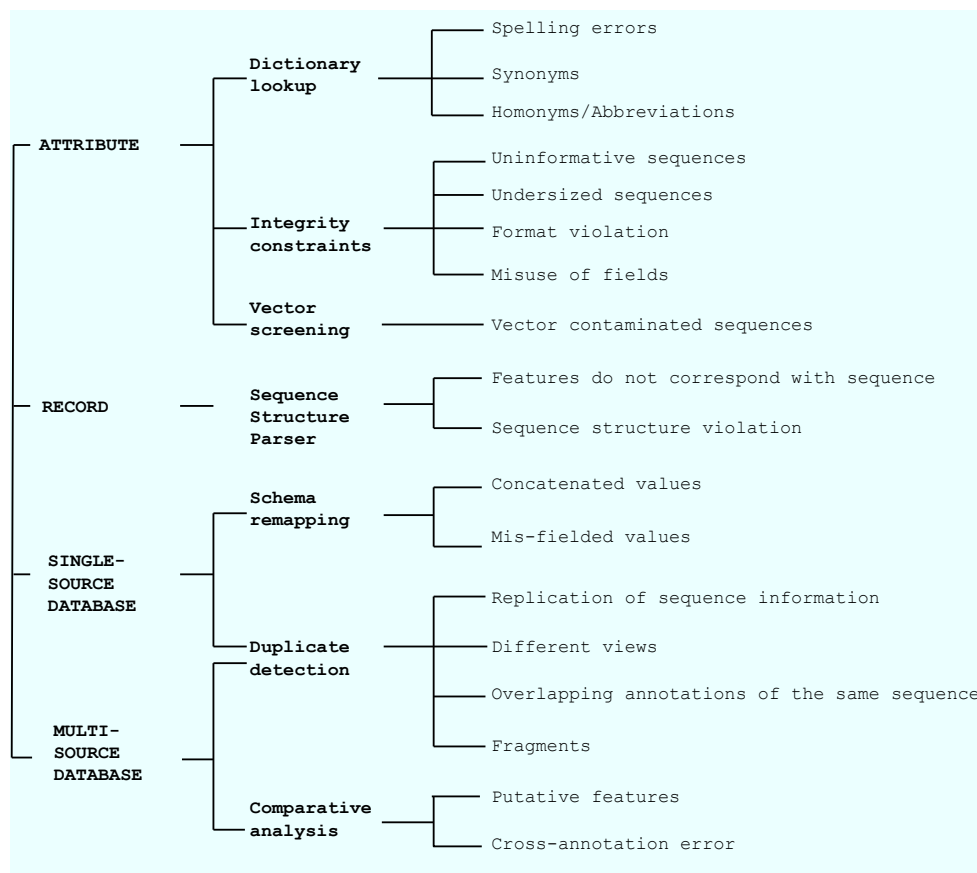


Figure 2: Classification of biological data issues by Koh et al. [28, 29].

3.3.1. Attribute-level data quality issues

Attribute level data issues are field values with uninformative, invalid, erroneous or ambiguous content. Koh et al. [29] observed four main types of attribute level data quality issues — invalid attribute values, ambiguous attribute values, dubious sequences, and contaminated sequences.

Invalid values: Header data issues result from the use of non-standard names, free-text entries, and from the diversity of database schema used in different databases. The header information is usually entered and provided by the person who submits the original record to the database (such as the description). Hence, the header information may contain spelling errors or invalid field values. Koh et al. [29] identified 569 possible misspelled words affecting up to 20,505 nucleotide records in Entrez. For example, “immunoglobulin” is misspelled as “immunogloblin” in the record with accession number AB122023⁹. Another example is the organism “brachydanio rerio” (zebrafish) which is misspelled “brachidanio rerio” as in the record with

⁹<http://www.ncbi.nlm.nih.gov/nuccore/AB122023>

1
2
3
4 accession numbers L25057¹⁰.

5
6
7 **Ambiguity:** As there is a lack of standardized naming conventions and controlled vocabulary use, vastly
8 different definitions may be used in database records to refer to the same sequence. The naming errors
9 include use of different names for the same sequence (synonym problem) or the same name for different
10 sequences (homonym problem) [53]. For example, the scorpion neurotoxin “BmK-X” precursor has many
11 possibly synonymous permutations. It is also known as “BmKX”, “BmK10”, “BmK-M10”, “Bmk M10”,
12 “Neurotoxin M10”, “Alpha-neurotoxin TX9”, and “BmKalphaTx9”¹¹.

13
14
15
16 Another type of error is the use of abbreviations, which may result in ambiguities. For example, the
17 abbreviation BMK stands for “Big Map Kinase”, “B-cell/myeloid kinase”, “bovine midkine”, as well as for
18 “Bradykinin-potentiating peptide”. GK is the abbreviation for both “Glycerol Kinase” and “Geko” gene of
19 *Drosophila melanogaster* (Fruit fly).

20
21
22 In free-text fields, a wrong piece of information can be entered as field value. For example, the description
23 of the sequence with the accession number AC254865¹² is “UNVERIFIED: BAC_10, complete sequence”.

24
25
26 **Dubious sequences:** The sequence is represented as a string of letters denoting the 20 amino acids in the
27 case of a protein sequence and the 4 nucleotides in the case of nucleotide sequence. Each base or residue is
28 thus limited to its alphabetical representation and “X” is used to denote an unknown residue, and “N” to
29 denote an unknown base. A base or residue which doesn’t correspond to its set of special letters is invalid and
30 can be caused by an erroneous data entry. For example, the sequence with the accession number AC000016¹³
31 contains 11% of unknown bases.

32
33
34
35
36 A sequence may also contain invalid symbols for nucleotides or amino acids, or it can be shorter than its
37 logical size. The length of protein sequences usually ranges from 6 to a few thousands residues. However,
38 Koh et al. [29] found 3,327 undersized protein sequences which are shorter than six residues in the public
39 databases using Entrez (as of Sep 2004), among which 1,887 contain only one residue.

40
41
42
43 **Contaminated sequences:** There are cases where a DNA sequence contains vectors used for cloning;
44 vector-contaminated sequences may be submitted to the database. Vectors are agents that carry DNA
45 fragments into a host cell. The vector sequences probe and bind the DNA fragments at the 5’ and 3’ sites.
46 The DNA fragment is then isolated from its vectors by cutting at the restriction enzyme sites. The existence
47 of vector-contaminated sequences was first reported in 1992; 0.23% of 20,000 eukaryotic entries were found
48 to be contaminated [20]. In 1999, [50] reported that up to 0.36% or 3,029 of the sequences in GenBank
49 contain contamination of the cloning vectors.

50
51
52
53
54
55 ¹⁰<http://www.ncbi.nlm.nih.gov/nuccore/L25057>

56 ¹¹<http://www.uniprot.org/uniprot/061705>

57 ¹²<http://www.ncbi.nlm.nih.gov/nuccore/AC254865>

58 ¹³https://www.ncbi.nlm.nih.gov/nuccore/AC000016.1?fmt_mask=65536

3.3.2. Record-level data quality issues

Conflicting information exists in the single record among two or more attributes — Koh et al. [28] call them the record-level data quality issues. Two types of record-level data quality issue are found in sequence records: sequence structure violations and inconsistent content with related references.

Sequence Structure Violations: It is known that a gene structure has a set of logical constraints, and any infringement of these constraints constitutes a possible feature issue. Such logical constraints include that introns and exons must be non-overlapping except in cases of alternative splicing. Koh et al. [29] observed that 12 out of 42,359 nucleotide sequences had overlapping intron/exon regions. For example, Syn7 gene of putative polyketide synthase in NCBI TPA record BN000507¹⁴ has overlapping intron 5 and exon 6. The rpb7+ RNA polymerase II subunit in GenBank record AF055916¹⁵ has overlapping exon 1 and exon 2.

Inconsistent with the literature: Usually, each record is associated with a list of publications by the authors of the sequence that discuss the data reported in the record. However, it is possible that a record is inconsistent with the information provided in the literature in general, and in the articles related to that record in particular.

For example, in the study of Dengue virus, Koh et al. [28] observed mis-annotations in Swiss-Prot record P27915¹⁶ and PIR record GNWVD3 [27]. The NS1/NS2A and NS4A/NS4B junctions given in these Dengue type 3 complete RNA sequences did not match the regions given in the reference of these records [38]. While manual checking of such inconsistencies by cross-referencing the database content with their corresponding literature is tedious, computational detection of discrepancies of the sequence annotations with its references is also non-trivial and may require complex text-mining solutions.

3.3.3. Single Database level data quality issues

Annotation errors: The features of a sequence are often directly submitted by the author of the sequence. The features can be derived experimentally or inferred. Computationally inferred features are usually based on sequence homologues and are derived using annotation tools. Hence, multiple database records of the same nucleotide or protein sequences may contain inconsistent or conflicting feature annotations. Koh et al. [29] refer to such data issues as cross-annotation errors. They identify possible causes of cross-annotation errors as: (i) different expert interpretations, (ii) mis-annotation of sequence function, and (iii) inference of features or annotation transfer based on best matches of low sequence similarity.

Annotation errors commonly result from mis-annotation or from data entry errors. In GenBank entries that contain splice site features in *Arabidopsis thaliana*, some 15% were found to have incorrect annotations [32]. Another study [24] found that 24% of the *Chlamydia trachomatis* sequences contained erroneous

¹⁴<https://www.ncbi.nlm.nih.gov/nuccore/BN000507>

¹⁵<https://www.ncbi.nlm.nih.gov/nuccore/AF055916>

¹⁶<http://www.uniprot.org/uniprot/P27915>

functional assignments. Another form of annotation error is caused by inaccurate inference of features from homologues.

Sequence duplication: Sequence duplication is also observed in sequence databases [12]. There are three types of redundancy: (1) the same sequence and annotations can be found in multiple records, (2) the same sequence but different annotations are found in multiple records, and (3) partially overlapping annotations of the same sequence exist in multiple databases which have different data views. For example, the records AAG39642¹⁷ and AAG39643¹⁸ contain identical sequences with exactly the same annotations.

3.3.4. Multiple Database level data quality issues

Due to the existence of heterogeneous database schema, massive data transformation is carried out in the databases during large-scale uploads or during data exchange. The transformation of data records from one schema to another may cause data integration problems, where data may be mapped to the wrong fields.

Finally, in this work we are interested in the detection of records that contain errors and inconsistencies through the analysis of the header section, and through a cross validation with the published literature. Thus, the model we built will not be able to detect errors related to the features or the sequence itself such as contaminations, undersized sequences, cross annotation errors, etc. Rather, it detects inconsistencies with the published literature.

3.4. Research problem statement

We propose to follow an IR approach, where a database record is regarded as a query, and its associated articles as the relevant documents. We use the term “query” to refer to a record, and the term “relevant documents” to refer to the set of its associated articles in its reference field.

We define the problem we study in this paper as follows. Given:

- a collection of documents $D = \langle d_1, d_2, \dots, d_n \rangle$
- a set of annotated queries $Q = \langle (q_1, y_1), (q_2, y_2), \dots, (q_m, y_m) \rangle$, where $y_m \in \{confident, suspicious\}$
- the set of relevant documents $D_R = \langle D_{R_1}, D_{R_2}, \dots, D_{R_m} \rangle$

we aim to retrieve and identify queries that are not consistent with their relevant documents or with the collection as whole, indicating that their description in the database record is incompatible with the information given in the corresponding publication, and which can therefore be flagged as “suspicious”. The resulting tool is expected to be used at curation time, and should send such “suspicious” records to curators for review.

¹⁷<https://www.ncbi.nlm.nih.gov/protein/AAG39642>

¹⁸<https://www.ncbi.nlm.nih.gov/protein/AAG39643>

4. Quality Factors

In this section, we introduce the features that we will consider as record quality indicators. We describe two kinds of features: record-based features, which mainly focus on the characteristics of the records; and IR-based features, which mainly focus on query quality predictors. Our approach here is to define a wide variety of features and then identify which of them are most valuable in classification.

4.1. Record-based features (9 features)

The characteristics of a record are strong indicators of its quality. We define several features that consider a record as a whole. Hence, we mainly rely on basic and intuitive quality factors, such as the record popularity, as well as building on recommendations given by the International Nucleotide Sequence Database Collaboration (INSDC) for the record structure.

Organism popularity (1 feature): Based on the intuition that organisms that have rarely been sequenced and deposited in a sequence database are more likely to have suspicious records, we consider the popularity of the main organism of a record as a quality feature. We define the popularity of an organism as the number of records that relate to that organism divided by the total number of records.

Record definition structure (3 features): The INSDC suggests that the record definition should have the following specific format:¹⁹ (i) it should start with the common name of the source organism; (ii) it gives the criteria by which this sequence is distinguished from the remainder of the source genome, such as the gene name and what it codes for, or the protein name and mRNA, or some description of the sequence's function (if the sequence is non-coding); (iii) if the sequence has a coding region, the description may be followed by a completeness qualifier, such as 'cds' (complete coding sequence). We define boolean attributes to indicate whether each of these rules is respected in a record or not.

Record popularity (1 feature): a popular sequence record is more expected to have been checked by other users, and hence be a confident record. Thus, we include the popularity as a quality feature and define it simply as the number of citations the record has.

Coding sequence (3 features): For a sequence with a coding region, the coding sequence (CDS) field in the features section of the records is one of the most important fields. Based on the feature table format designed jointly by GenBank, the EMBL Nucleotide Sequence Data Library, and the DNA Data Bank of Japan,²⁰ the CDS field should specify: (i) the region of nucleotides that corresponds with the sequence of amino acids in a protein (location including start and stop codons), (ii) the gene name, and (iii) the product/protein name. For each CDS field of a record, we check its validity (i) by ensuring that the CDS region is within the sequence range, (ii) by ensuring that the gene name is valid and is given into the

¹⁹<ftp://ftp.ncbi.nih.gov/GenBank/gbrel.txt>

²⁰<http://www.insdc.org/documents/feature-table>

1
2
3
4 annotations of Gene Ontology (GO) ([54]), and (iii) by ensuring that the protein exists. Hence, we define
5 three quality attributes for (i), (ii), and (iii) based on an aggregation of all CDS fields of each record.
6

7
8 **Definition length (1 feature):** the length of the definition may indicate the degree of precision given to
9 describe a record. Hence, we include the length as a quality factor, and we define it as the total number of
10 terms.
11

12 13 14 4.2. IR-based features (203 features)

15 To find indicators or features to represent the quality of each query (record), we draw on the large body
16 of previous work on query quality prediction ([16, 22, 33]). While some of these features such as Overlap
17 Similarity are stand-alone, other features such as Average TF are derived from term level statistics ([33]).
18 These include predicting the quality of queries using either pre-retrieval indicators like Query Scope, that
19 is, they are calculated for a query as a whole, or post-retrieval indicators like Query Clarity, that is, they
20 involve performing an initial retrieval and hence are more expensive to compute. We describe the set of
21 query quality predictors we used. As stated previously, to compute these IR-based features, we consider the
22 record *definition* field as the query.
23
24
25
26
27

28
29 **Query clarity (QC) (18 features):** Developed by [16], this post-retrieval factor is the Kullback-Leibler
30 divergence of the query model from the collection model. QC is computed as:
31

$$32 \quad QC = \sum_{w \in q} P(w|q) \times \log_2 \frac{P(w|q)}{P_C(w)} \quad (1)$$

33 where $P(w|q)$ is the probability of the occurrence of the word w in the query model, and $P_C(w)$ is the
34 probability of the occurrence of w in the collection. The query model is estimated from the top- k ranked
35 documents retrieved after an initial run of the original query. We computed different QC scores based on
36 different configuration options $k \in \{1, 5, 10, 20, 50, 100\} \times w \in \{\text{title, abstract, body}\}$.
37
38
39
40
41

42
43 **Simplified clarity score (SCS) (3 features):** To avoid the expensive computation of query clarity, [21]
44 proposed simplified clarity score as a comparable pre-retrieval performance factor. It is calculated as:
45
46

$$47 \quad SCS = \sum_{w \in q} P_{ml}(w|q) \times \log_2 \frac{P_{ml}(w|q)}{P_C(w)} \quad (2)$$

48 where $P_{ml}(w|q)$ is the probability of the occurrence of the word w in the query. We also computed SCS
49 based on different configuration options of $word \in \{\text{title, abstract, body}\}$.
50
51
52
53

54
55 **Relevant-documents clarity score (RDCS) (3 features):** We also propose to compute the clarity score
56 based on a query model estimated from the relevant documents themselves, while considering separately three
57 different fields of the relevant documents $\{\text{title, abstract, body}\}$.
58
59

60 **IDF-based features (24 features):** We calculate the IDF of each query term w as:
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

$$IDF_w = \log\left(1 + \frac{N}{N_w}\right) \quad (3)$$

where N_w is the document frequency of w while considering separately three different fields {title, abstract, body}, and N is the number of documents in the collection. For each query we calculate the sum, standard deviation, minimum, maximum, arithmetic mean, geometric mean, harmonic mean, and coefficient of variation of the IDF's of constituent terms.

TF-based features (24 features): We calculate the TF of each query term w in a relevant document d as:

$$TF_w = \log(1 + f_{w,d}) \quad (4)$$

where $f_{w,d}$ is the number of time w occurs in d while considering separately three different fields {title, abstract, body}. For each TF value of each term, we calculate aggregate values similar to those for IDF as quality factors.

Similarity of collection–query (SCQ) (24 features): Proposed by [59], this query quality factor is based on the hypothesis that queries that have higher similarity to the collection as a whole will be of higher quality. For each term w in the query, SCQ is computed as:

$$SCQ_w = (1 + \ln(n(w)) \times \ln\left(1 + \frac{N}{N_w}\right)) \quad (5)$$

where $n(w)$ is the frequency of the term w in the collection while also considering separately three different fields {title, abstract, body}. Based on the SCQ values of each term, we also calculated aggregate values similar to those for IDF.

Inverse collection term frequency features (ICTF) (24 features): The inverse collection term frequency of a term w is defined as:

$$ICTF_w = \log\left(1 + \frac{T}{n(w)}\right) \quad (6)$$

where $n(w)$ is the frequency of the term w in the collection and T is the number of term occurrences in the collection while considering separately three different fields {title, abstract, body}. Using the ICTF values, we calculate aggregate statistics similar to those for IDF.

Query scope (QS) (4 features): Query scope ([21]) is a measure of the size of the retrieved document set relative to the size of the collection. We can expect that high values of query scope are predictive of poor-quality queries as they retrieve far too many documents. QS is computed as:

$$QS = \log\left(1 + \frac{N}{n_q}\right) \quad (7)$$

1
2
3
4 where n_q is the number of documents that match the query terms while considering separately four different
5 fields {title, abstract, body, all document}.

6
7
8 **Similarity of relevant documents–query features (48 features):** Based on the fact that a high
9 similarity value between a query and its relevant documents reflects a high query quality, we include several
10 information-theoretic, statistical, and practical similarity measures as quality indicators. These similarity
11 measures are: matching, overlap, Jaccard, Dice, cosine, mutual information (MI), and Okapi BM25. We also
12 used various IR similarity ranking functions including: the sum of TFIDF scores (SumTFIDF), the Lucene
13 vector-space model score (LuceneVSM),²¹ the BM25 score ([45]), language model scores based on (i) the
14 Jelinek-Mercer smoothing (LMJelinekMercer) ([58]) and on (ii) a Bayesian smoothing using Dirichlet priors
15 (LMDirichlet) ([58]), and an information-based score (IBSimilarity) ([13]). These similarities are computed
16 while considering separately four different fields {title, abstract, body, all document}.

17
18
19
20
21
22
23 **Retrieval performance score (RPS) (28 features):** Based on the relevance paradigm of IR, we assume
24 that a good quality record should rank its corresponding articles highly. Thus, we use the reciprocal rank
25 evaluation measure to define the RPS as follows:

$$RR = \frac{1}{rank_i} \quad (8)$$

26
27
28
29
30
31
32 where $rank_i$ is the rank of the first relevant document in the retrieved list of documents that matches the
33 query q_i returned by the system. We have also considered query expansion using the following terms related
34 to the organism: (i) scientific name, (ii) common names, (iii) synonyms, (iv) abbreviations, (v) misnomers,
35 and (vi) misspellings. These terms are extracted from the NCBI Taxonomy,²² which is a curated classification
36 and nomenclature for all of the organisms in the public sequence databases. The basic intuition is that: if
37 (i), (ii), (iii), and (iv) improve the retrieval performance, there is a mismatch between the record and its
38 corresponding article, and thus, the record may be of low quality. Also, if (v) and (vi) improve the retrieval
39 performance, the article is clearly reporting the record using incorrect terms, and thus, the record is probably
40 of low quality. Here we consider querying separately four different fields {title, abstract, body, all document}.

41
42
43
44
45
46 **Citation of the main organism in the article (3 features):** The citation of the main organism in
47 the relevant document is an important quality factor since the article is supposed to report the content of
48 the record. Hence, we include this quality factor a binary feature while considering separately four different
49 fields {title, abstract, body}.

50
51
52
53 In total, we have defined 9 record-based features and 203 IR-based features, for a total of 212 features
54 that characterize the quality records.

55
56
57 ²¹https://lucene.apache.org/core/6_1_0/core/org/apache/lucene/search/similarities/TFIDFSimilarity.html

58 ²²<http://www.ncbi.nlm.nih.gov/taxonomy>

5. Experimental Setup

In this section, we describe the dataset we have constructed from publicly available resources, and then introduce the learning algorithm we used to classify records as “*confident*” or “*suspicious*”.

5.1. Data description

Articles: We used the PubMed Central Open Access collection²³ (OA), which is a free full-text archive of biomedical and life sciences journal literature at the U.S. National Institutes of Health’s National Library of Medicine. The release of PMC OA we used contains roughly 1.13 million articles, which are provided in an XML format with specific fields corresponding to each section or subsection in the article. We used the Lucene IR System²⁴ to index the collection, with the default settings for stemming and English stop-word removal. We defined a list of biomedical keywords, which should not be stemmed or considered as stop-words, such as the protein names “*THE*” and “*Is*”. Each section of an article (title, abstract, body) is indexed separately, so that different sections can be used and queried separately to compute the quality features.

Sequences: We work with the GenBank nucleotide database, but limit the sequence database records we work with to those that are cited by the PMC OA article collection. Specifically, we used a regular expression to extract GenBank accession numbers mentioned in the PMC OA articles, thereby identifying literature that refers to at least one GenBank identifier. This resulted in a list of 733,779 putative accession numbers. Of these, 494,142 were valid GenBank nucleotide records that we were able to download using the e-utilities API ([47]).²⁵ Among the valid records, only 162,913 records also cite the corresponding articles (as determined by matching their titles). This process gave us a list of 162,913 pairs of record accession numbers and PMC article identifiers, which cite each other. Note that for the 331,229 records that we have put aside, each record cites an article; however, we do not have access to all articles through PMC OA.

Each record in this dataset was labelled as “*alive*” or “*dead*”, an attribute that we obtained using the eSummary API ([47]). Note that the records that are reported as “*dead*” are explicitly labelled as such in GenBank, whereas records that are “*alive*” are implicitly labelled by not being dead. In the classification task, we consider dead records to be “*suspicious*” and all other records to be “*confident*” in our labelling. We acknowledge that an “*alive*” record does not necessarily indicate that it is of good quality. However, we made this assumption motivated by the fact that, overall, the records are of good quality, whereas only a small fraction of the data may be faulty. This has been observed and reported in [6], where the authors carried out a biocuration task by manually analyzing records that are alive. The authors have found that among 100 alive records randomly selected in the dataset, roughly 5% have been reported as being faulty by the database curator. Hence, we believe that even if this very small fraction of records that are faulty

²³<http://www.ncbi.nlm.nih.gov/pmc/tools/openftlist/> The version used was downloaded on October 2015.

²⁴<http://lucene.apache.org/>

²⁵The sequences were downloaded on October 2015.

1
2
3
4 and are included in our training dataset as of "good quality", they would have an insignificant impact on
5 the learned model.
6

7 Among the 162,913 records for which the relevant articles are in the PMC OA dataset, 162,486 are alive
8 and only 427 are dead. Hence, our dataset is skewed toward negative examples (alive records) with only a
9 few positive examples (dead records).
10
11

12 **Organism taxonomy:** To gather more information about the record organisms, such as the list of syn-
13 onyms, acronyms and common names, we used the NCBI Taxonomy database. This is a curated classification
14 of all of the organisms in the public sequence databases. It represents about 10% of the described species of
15 life on the planet.
16
17
18

19 20 5.2. Anomaly detection algorithm

21
22 Given as input a set of quality indicators for each record, our goal is to combine these inputs to produce a
23 value indicating whether the record is "confident" or "suspicious". To accomplish this, we used the Support
24 Vector Machines (SVM) classification algorithm ([14]), which is one of the most widely-used and effective
25 classification algorithm.
26
27

28 Each record m is represented by its vector of k quality indicators $x_m = [x_{m1}, x_{m2}, \dots, x_{mk}]$ and its asso-
29 ciated label $y_m \in \{\text{confident}, \text{suspicious}\}$. We used the SVM implementation available in the LibSVM ([9])
30 package. Both Linear and RBF kernels were considered in our experiments. The regularization parameter C
31 (the trade-off between training error and margin), the gamma parameter of the RBF kernel, and the penalty
32 parameter w_i that penalizes negative examples due to the skewed nature of the dataset were selected from
33 a search within the discrete sets $\{10^{-5}, 10^{-3}, \dots, 10^{13}, 10^{15}\}$, $\{10^{-15}, 10^{-13}, \dots, 10^1, 10^3\}$, and $\{10, 20, \dots,$
34 $50, 100, 200\}$ respectively, using 10 fold cross validation.
35
36
37
38

39 Although the differences were not substantial, experiments with the best RBF kernel parameters per-
40 formed slightly better than the best linear kernel parameters for the majority of the validation experiments.
41 Unless otherwise noted, all presented results were obtained using an RBF kernel, with C set to 10^{-3} , gamma
42 set to 10^{-3} , and w_i set to 100.
43
44
45
46

47 6. Experimental Evaluation

48
49 We now report and discuss the main results of the experimental evaluation, considering both the effec-
50 tiveness of the method and our interpretation of which features are valuable in classification.
51
52

53 6.1. Feature analysis

54
55 To explore the relationship between features and the record quality labels, we undertook a feature analysis
56 task. A general method for measuring the amount of information that a feature x_k provides w.r.t. predicting
57 a class label y ("confident" or "suspicious") is to calculate its mutual information (MI) $I(x_k, y)$ or Pearson's
58 chi-squared test $\chi^2(x_k, y)$. In Table 1, we present the list of ten top-ranked features using these two metrics.
59
60
61
62

Table 1: Ranking of the most important features using two different metrics.

	Mutal Information		Pearson's chi-squared test	
Rank	Field	Feature	Field	Feature
1	title	LMDirichlet Score	title	LMDirichlet Score
2	abstract	SumTFIDF score	abstract	SumTFIDF Score
3	abstract	LMDirichlet Score	abstract	RDCS
4	abstract	BM25 Score	abstract	LMDirichlet Score
5	abstract	RDCS	abstract	BM25 Score
6	abstract	IBSimilarity Score	abstract	IBSimilarity Score
7	title	SumTFIDF Score	title	BM25 Score
8	title	BM25 Score	title	IBSimilarity Score
9	title	IBSimilarity Score	title	SumTFIDF Score
10	x	<i>Popularity Organism</i>	<i>abstract</i>	<i>LMJelinekMercer Score</i>

The two lists are roughly similar except for the tenth line, where MI introduces a record-based feature. These two lists led us to make the following observations:

1. Features based on the similarity between the relevant documents and the record are the most informative (8/10 for mutual information and 9/10 for chi-squared). This confirms that a good and a confident record is highly discussed in its associated articles.
2. IR similarity ranking functions are the most informative features. They take into account the information carried in both the query and the documents, in contrast to statistical similarity measures.
3. For both rankings, the top feature is a language-model similarity score. It computes the similarity between the record definition and the title of the relevant document, using Bayesian smoothing with Dirichlet priors. This shows that a confident record is one in which the description has a high probability of having been generated from the title of its associated document.
4. Top features are mainly based on short and medium document fields (that is, title and abstract). This reflects the fact that confident records can be expected to be referenced and discussed earlier in the article.
5. Almost all top features are IR-based features. Only one record-based feature appears in the two rankings (popularity of the organism in the MI ranking). This confirms our initial assumption that the literature is a strong resource for checking the quality of a record.

6.2. Classification performance

The effectiveness of our method (denoted SVM LBF+RBF, for SVM with Literature-based Features and Record-based Features), is summarized in Table 2, broken down by the two classes in the data. The last

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 2: The results of the classification accuracy.

Algorithm	Confident records						Suspicious records						Harmonic Mean of F1-Scores	Improvement rate (%)
	Precision	Recall	F1-Score	TN	FN	FP	Precision	Recall	F1-Score	TP	FP			
SVM LBF+RBF ^a	0.998	0.993	0.995	161,380	300	1,106	0.103	0.299	0.154	127	1,106	0.267	-	
SVM RBF ^b	0.997	0.982	0.997	162,195	412	291	0.049	0.035	0.041	15	291	0.078	242.31%	
SVM LBF ^c	0.998	0.989	0.993	160,746	312	1,740	0.062	0.269	0.101	115	1,740	0.183	45.90%	
Majority class	1.000	1.000	1.000	162,486	427	0	0.000	0.000	0.000	0	0	0.000	+∞	
Random classification 1	0.997	0.499	0.666	57,696	217	104,790	0.002	0.494	0.005	210	104,790	0.009	2866%	
Random classification 2	0.997	0.997	0.997	161,987	426	499	0.002	0.002	0.002	1	499	0.003	8800%	
RPS-based classification	0.997	0.927	0.960	150,285	378	12,201	0.004	0.114	0.007	49	12,201	0.013	1953%	

^aSVM LBF+RBF: using both Record-Based Features and Literature-Based Features. ^bSVM RBF: using only record-based features. ^cSVM LBF: using only IR-based features.

column of the table shows the harmonic mean between the F1-Scores of the two classes. We provide a comparison with two variants of this method and four baseline methods:

- SVM RBF: SVM classifier trained with only record-based features.
- SVM LBF: SVM classifier trained with only IR-based features.
- Majority class: a naïve approach that simply predicts the most common class in the data (that is, classify everything as “*confident*”).
- Random classification 1: which classifies a record as “*confident*” or “*suspicious*” with a 50% probability. The classification was performed 10 times independently over the full dataset; we report average results.
- Random classification 2: classifies a record as “*confident*” or “*suspicious*” with a 99.73% probability of being classified as confident. This value reflects the natural distribution of the source data, since 99.73% of the records are ‘alive’. The classification was also performed 10 times independently over the full dataset; we report average results.
- RPS-based classification: classifies a record as “*confident*” or “*suspicious*” based on a fixed threshold (0.05) for the RPS value.

Our method (SVM LBF+RBF) shows a statistically significant improvement over the best baseline (relatively, 45.90% over SVM LBF). Second, the results confirm again our initial assumption that the literature is a strong support to assess the quality of records in sequence databases. Third, by comparing SVM RBF and SVM LBF, we conclude that the associated literature provides better evaluation of quality than can be obtained by examining only the records. This has also been shown in the previous section through the feature analysis. Due to the skewed nature of the dataset, all algorithms tend to classify the records as “*confident*”, which results in high precision and recall values for this class, but these results are not very informative. It is more meaningful to consider the performance over the class of “*suspicious*” records.

The instance-level TP and FP values given in Table 2 for each method illustrate how many FP records would need to be reviewed in order to find the small number of TPs in each case. This allows to compare the curator task difficulty for each method; the curator would need to examine everything retrieved (all Positives) and then make a decision as to whether it was a correct retrieval of a suspicious record (TP) or a perfectly valid record (FP from the perspective of “*suspicious*”). The table shows how much work needs to be done in each scenario, and demonstrates that our method considerably reduces the curation workload. In particular, to get reasonable recall with the random approach, the curator needs to review more than 100,000 records. Indeed, the curator would need to review 85x as many records with a random approach as compared to our method, for a gain of only 83 TP records (respectively: 1,233 vs. 105,000 positive records; 127 vs. 210 true positives identified). This highlights the substantial amount of effort saved using the approach we are proposing.

1
2
3
4 Table 2 shows relatively low values for precision and recall compared with some other machine learning
5 problems. However, first, the dataset is highly imbalanced, with far more records labelled as “*confident*” than
6 labelled as “*suspicious*”. Second, the records that are labelled as “*suspicious*” have been explicitly labelled,
7 whereas records that are “*alive*” are implicitly labelled and assumed (perhaps wrongly) to be “*confident*”;
8 an alive record may be a low-quality record that was missed in error. Hence, training a learning algorithm on
9 unlabelled data leads to poor effectiveness, particularly when the minority class is the most likely class to be
10 missed. Probably some of the records which have been classified as “*suspicious*” by our learning algorithm,
11 but were labelled in the data set as “*confident*”, are in fact problematic. We discuss below typical examples
12 of records which have been incorrectly classified by our method (false positives and false negatives). The
13 profile of each example is given in Table 3, using the top 7 features obtained in the feature analysis in Section
14 6.1.
15
16
17
18
19
20
21

22 **Example 1.** (False Positive) The record with accession number FJ824848²⁶ has been classified as “*suspi-*
23 *cious*” by the algorithm, and presents the typical profile of a suspicious record as given in Table 3. This
24 record presents the complete sequence of the cloning vector pDMK3.²⁷ First, the record definition does not
25 give much information. Second, this record presents an organism for which there are relatively few other
26 records. Third, the content of the record is mentioned neither in the title of the article nor in the abstract
27 of the article with PMC identifier PMC2675058²⁸. By examining the content of the article, we have noted
28 that this cloning vector is mentioned as pDMK2. This case leads us to make two inferences: either the title
29 of the record is incorrect, or the article uses an incorrect term to refer to this cloning vector. Consequently,
30 we contacted the corresponding author of PMC2675058, Anders Sjöstedt, who acknowledged the error by
31 saying: “For practical purposes this doesn’t matter since the two vectors are identical with the exception
32 of two additional restriction sites in pDMK3. We should have stated pDMK3 in the paper so it can be de-
33 noted as a typographical error.” (personal communication, Anders Sjöstedt). However, since both pDMK2
34 and pDMK3 exist, we believe that this error cannot be considered as typographical error, but rather as a
35 confusion and inconsistency between that record and its associated article.
36
37
38
39
40
41
42
43
44
45

46 **Example 2.** (False Positive) The record with accession number CP006742²⁹ has also been classified as
47 “*suspicious*” by the algorithm, and also presents the typical profile of a suspicious record as given in Table
48 3. This record presents the complete genome of the *Bacillus anthracis* organism. In fact, according to the
49 article that reports the content of that record (accession number PMC3923885³⁰), this record is supposed
50 to present the chromosome Cow1. However, this chromosome is not mentioned in the record. We have
51
52
53
54

55 ²⁶<http://www.ncbi.nlm.nih.gov/nuccore/FJ824848>

56 ²⁷A cloning vector is a small piece of DNA, taken from a virus, a plasmid, or the cell of a higher organism, that can be stably
57 maintained in an organism, and into which a foreign DNA fragment can be inserted for cloning purposes.

58 ²⁸<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2675058/>

59 ²⁹<http://www.ncbi.nlm.nih.gov/nuccore/CP006742>

60 ³⁰<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3923885/>

Table 3: Example profiles.

		Example 1	Example 2	Example 3	Example 4
	Accession	FJ824848	CP006742	BK008760	DW407270
	PMC	PMC2675058	PMC3923885	PMC4233938	PMC1621083
	Type	FN	FN	FN	FP
1	Tit. LMDir.	0.00	0.00	1.4194	18.78
2	Abs. SumTFIDF	0.00	0.00	2.25	27.04
3	Abs. LMDir.	0.00	0.00	3.3450	30.26
4	Abs. BM25	0.00	0.00	2.80	65.00
5	Abs. RDCS	0.00	0.0302	0.034	0.68
6	Pop. Organism	0.00000613	0.00	2.4828	0.00086
7	RPS	0.00	0.0019	0.20	1.00

contacted the corresponding author, and he answered as follows: “After some reasoning, we decided to call the isolates cow 1,2,3.. etc in the paper to increase readability. Their names in our strain collection is however different. As far as I can see CP006742 is the B. anthracis chromosome sequence that we submitted, although I haven’t checked every basepair...” (personal communication, Bo Segerman). Hence, we believe that this little inconsistency between the record and its associated article has prompted the algorithm to classify this record as “*suspicious*”.

Example 3. (False Positive) The record with accession number BK008760³¹ has also been classified as “*suspicious*” by the algorithm. This record presents relatively low values in its profile compared to the false negative example presented as given in Table 3. According to the article that reports the content of that record (PMC identifier PMC4233938³²), this article is supposed to show two genes, Atg8a and Atg8b. However, the record is only showing the coding sequence of the gene Atg8a. We also have contacted the paper’s authors who acknowledged the error by saying: “Thanks for your important notification. Indeed this is an entry error. The error was corrected in a Corrigendum that was published soon after the original paper” [48] (personal communication, Assaf Vardi). Although the error is again in the research article itself, this shows another record-literature inconsistency example, which illustrates a false positive that contributed to the low precision-recall values obtained.

Example 4. (False Negative) In contrast, the record with the accession number DW407270³³ has been classified “*confident*” by our method, while its current status is “*suspicious*”. The record references a popular organism and is correctly discussed in its associated article. We did not discover any mismatch or conflict

³¹<http://www.ncbi.nlm.nih.gov/nuccore/BK008760>

³²<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4233938/>

³³<http://www.ncbi.nlm.nih.gov/nuccore/DW407270>

1
2
3
4 between this record and its associated article (PMC identifier PMC1621083³⁴). In fact, this record has been
5 removed because the underlying biological material suffered from bacterial contamination with *Pseudomonas*
6 *fluorescens*, an issue that could only be resolved upon consideration of the similarity of the record sequence
7 to other sequences via a BLAST search. This example leads us to identify an important distinction — a
8 record can be coherent from one perspective (literature consistency) while being inconsistent from another
9 (biological content). It highlights the need to consider data quality from more than one perspective, and
10 demonstrates that literature alone is insufficient to detect all suspicious records. It further explains the
11 limitations on performance of our method.
12
13
14
15
16

17 The examples discussed above in this section are not necessarily faulty records. We have rather presented
18 these examples to provide a brief justification for why we obtained low precision-recall values. This motivates
19 us to build a manually curated dataset in order to remove any ambiguous and noisy records that may lead
20 to the build of a biased model, and to develop features that will capture further aspects of sequence record
21 quality. This will form part of our future work.
22
23
24
25
26

27 7. Conclusions and Future Work

28
29 In this paper we have introduced a list of factors that correlate with the quality of a record. We used
30 these quality indicators to train an anomaly detection algorithm based on supervised learning to classify
31 records as “*confident*” or “*suspicious*”. We then performed a complete analysis on the full PubMed Central
32 collection. The main outcome of this work is evidence that, in addition to the sequence itself, the literature
33 is a valuable external resource that can be used to assess the quality of a database record.
34
35
36

37 Despite the fact that our method significantly outperforms the suggested baselines, we obtained some-
38 what low effectiveness scores, compared to some other common machine learning problems. Therefore, we
39 undertook a feature analysis and a failure analysis, examining specific cases that indicate causes of this
40 performance — in particular, identifying that the ground truth may contain errors as well as recognising
41 that literature alone is insufficient to represent the full spectrum of data quality issues.
42
43
44

45 This work is to the best of our knowledge the first use of the literature as a tool for addressing the
46 data quality problem in biomedical sequence databases. We have shown that the approach can identify
47 problematic records with enough accuracy to be of value to curators, potentially reducing the effort required
48 to remove low-quality records by nearly two orders of magnitude.
49
50

51 Our current dataset relies on data from GenBank to obtain labels, and the negative labels are only derived
52 implicitly. This suggests two directions for future work. First, it would be desirable to construct a manually
53 curated dataset explicitly for development of automated quality analysis techniques. Second, there is a need
54 for new unsupervised learning methods for anomaly detection. It may be, for example, that good and bad
55 records have distinct distributions of attribute values, so that methods such as k-nearest neighbour or local
56
57
58
59

60 ³⁴<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1621083/>
61
62

1
2
3
4 outlier factor ([8]) could be applied. Having established that automated literature analysis can be applied
5 in practice to this task, the challenge now is to improve performance and further reduce the effort needed to
6 clean databases. We also expect that leveraging external textual information to support data cleaning will
7 have broader application in other database contexts.
8
9

10 11 12 **References**

- 13
14 [1] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool.
15 *Journal of Molecular Biology*, 215(3):403 – 410, 1990.
16
17 [2] B. P. Anton, S. Kasif, R. J. Roberts, and M. Steffen. Objective: biochemical function. *Frontiers in*
18 *genetics*, 5:210, 2014.
19
20 [3] D. Apiletti, G. Bruno, E. Ficarra, and E. Baralis. Data cleaning and semantic improvement in biological
21 databases. *Journal of Integrative Bioinformatics*, 3(2):40, 2006.
22
23 [4] F. B. Bastian, M. C. Chibucos, P. Gaudet, M. Giglio, G. L. Holliday, H. Huang, S. E. Lewis, A. Niknejad,
24 S. Orchard, S. Poux, et al. The confidence information ontology: a step towards a standard for asserting
25 confidence in annotations. *Database*, 2015:bav043, 2015.
26
27 [5] M. J. Bell, M. Collison, and P. Lord. Can inferred provenance and its visualisation be used to detect
28 erroneous annotation? a case study using uniprotkb. *PloS one*, 8(10):e75541, 2013.
29
30 [6] M. R. Bouadjenek, K. Verspoor, and J. Zobel. Literature consistency of bioinformatics sequence
31 databases is effective for assessing record quality. Technical report, The University of Melbourne, 2016.
32 <http://people.eng.unimelb.edu.au/mbouadjenek/papers/Biocuration2017.pdf>.
33
34 [7] S. E. Brenner. Errors in genome annotation. *Trends in Genetics*, 15(4):132 – 133, 1999.
35
36 [8] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. LOF: Identifying density-based local outliers. In
37 *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, SIGMOD
38 '00, pages 93–104, New York, NY, USA, 2000. ACM.
39
40 [9] C.-C. Chang and C.-J. Lin. Libsvm: A library for support vector machines. *ACM Trans. Intell. Syst.*
41 *Technol.*, 2(3):27:1–27:27, May 2011.
42
43 [10] S. Chellamuthu and M. Punithavalli. Detecting redundancy in biological databases? An efficient ap-
44 proach. *Global Journal of Computer Science and Technology*, 9(5):141–145, Sept. 2009.
45
46 [11] I.-M. A. Chen, V. M. Markowitz, K. Chu, I. Anderson, K. Mavromatis, N. C. Kyrpides, and N. N.
47 Ivanova. Improving microbial genome annotations in an integrated database context. *PLoS One*,
48 8(2):e54859, 2013.
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 [12] Q. Chen, J. Zobel, and K. Verspoor. Evaluation of a machine learning duplicate detection method for
5 bioinformatics databases. In *DTMBIO*, pages 4–12, New York, NY, USA, 2015. ACM.
6
7
8 [13] S. Clinchant and E. Gaussier. Information-based models for ad hoc ir. In *Proceedings of the 33rd*
9 *International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR
10 '10, pages 234–241, New York, NY, USA, 2010. ACM.
11
12
13 [14] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
14
15 [15] J. Crappé, E. Ndah, A. Koch, S. Steyaert, D. Gawron, S. De Keulenaer, E. De Meester, T. De Meyer,
16 W. Van Criekinge, P. Van Damme, et al. Proteoformer: deep proteome coverage through ribosome
17 profiling and ms integration. *Nucleic acids research*, page gku1283, 2014.
18
19
20 [16] S. Cronen-Townsend, Y. Zhou, and W. B. Croft. Predicting query performance. In *Proceedings of*
21 *the 25th Annual International ACM SIGIR Conference on Research and Development in Information*
22 *Retrieval*, SIGIR '02, pages 299–306, New York, NY, USA, 2002. ACM.
23
24
25 [17] M. I. Dunitz, J. M. Lang, G. Jospin, A. E. Darling, J. A. Eisen, and D. A. Coil. Swabs to genomes: a
26 comprehensive workflow. *PeerJ*, 3:e960, 2015.
27
28
29 [18] W. R. Gilks, B. Audit, D. De Angelis, S. Tsoka, and C. A. Ouzounis. Modeling the percolation of
30 annotation errors in a database of protein sequences. *Bioinformatics*, 18(12):1641–1649, 2002.
31
32
33 [19] J. Gillis and P. Pavlidis. Characterizing the state of the art in the computational assignment of gene
34 function: lessons from the first critical assessment of functional annotation (cafa). *BMC bioinformatics*,
35 14(3):1, 2013.
36
37
38 [20] R. Guigo, P. Agarwal, J. F. Abril, M. Burset, and J. W. Fickett. An assessment of gene prediction
39 accuracy in large dna sequences. *Genome Research*, 10(10):1631–1642, 2000.
40
41
42 [21] B. He and I. Ounis. Inferring query performance using pre-retrieval predictors. In *SPIRE*, pages 43–54.
43 Springer Berlin Heidelberg, 2004.
44
45
46 [22] B. He and I. Ounis. Query performance prediction. *Information Systems*, 31(7):585 – 594, 2006. (1)
47 SPIRE 2004(2) Multimedia Databases.
48
49
50 [23] L. Holm and C. Sander. Removing near-neighbour redundancy from large protein sequence collections.
51 *Bioinformatics*, 14(5):423–429, 1998.
52
53
54 [24] I. Iliopoulos, S. Tsoka, M. A. Andrade, A. J. Enright, M. Carroll, P. Poulet, V. Promponas, T. Liakopou-
55 los, G. Palaios, C. Pasquier, S. Hamodrakas, J. Tamames, A. T. Yagnik, A. Tramontano, D. Devos,
56 C. Blaschke, A. Valencia, D. Brett, D. Martin, C. Leroy, I. Rigoutsos, C. Sander, and C. A. Ouzounis.
57 Evaluation of annotation strategies using an entire genome sequence. *Bioinformatics*, 19(6), 2003.
58
59
60
61
62
63
64
65

- 1
2
3
4 [25] I. Kahanda, C. S. Funk, F. Ullah, K. M. Verspoor, and A. Ben-Hur. A close look at protein function
5 prediction evaluation protocols. *GigaScience*, 4(1):1, 2015.
6
7
8 [26] N. Kaplan and M. Linial. Automatic detection of false annotations via binary property clustering. *BMC*
9 *Bioinformatics*, 6(1):1–8, 2005.
10
11
12 [27] A. M. Khan, A. T. Heiny, K. X. Lee, K. N. Srinivasan, T. W. Tan, J. T. August, and V. Brusic. Large-
13 scale analysis of antigenic diversity of t-cell epitopes in dengue virus. *BMC Bioinformatics*, 7(S-5),
14 2006.
15
16
17 [28] J. L. Y. Koh. Correlation-based methods for biological data cleaning. Master’s thesis, School of Com-
18 puting National University of Singapore, 2007.
19
20
21 [29] J. L. Y. Koh, M. L. Lee, and V. Brusic. A classification of biological data artifacts. In *Workshop on*
22 *Database Issues in Biological Databases*, pages 53–57, 2005.
23
24
25 [30] J. L. Y. Koh, M. L. Lee, A. M. Khan, P. T. J. Tan, and V. Brusic. Duplicate detection in biologi-
26 cal data using association rule mining. In *European Workshop on Data Mining and Text Mining in*
27 *Bioinformatics*, pages 35–41, 2004.
28
29
30
31 [31] E. Koonin and M. Galperin. Sequence-evolution-function: Computational approaches. *Comparative*
32 *Genomics*, 2002.
33
34
35 [32] P. G. Korning, S. M. Hebsgaard, P. Rouzé, and S. Brunak. Cleaning the genbank arabidopsis thaliana
36 data set. *Nucleic Acids Research*, 24(2):316–320, 1996.
37
38
39 [33] G. Kumaran and V. R. Carvalho. Reducing long queries using query quality predictors. SIGIR ’09,
40 pages 564–571, New York, NY, USA, 2009. ACM.
41
42
43 [34] D. Lee, O. Redfern, and C. Orengo. Predicting protein function from sequence and structure. *Nature*
44 *Reviews Molecular Cell Biology*, 8(12):995–1005, 2007.
45
46
47 [35] W. Li and A. Godzik. Cd-hit: a fast program for clustering and comparing large sets of protein or
48 nucleotide sequences. *Bioinformatics*, 22(13):1658–1659, 2006.
49
50
51 [36] F. Mao, Z. Su, V. Olman, P. Dam, Z. Liu, and Y. Xu. Mapping of orthologous genes in the context of
52 biological pathways: An application of integer programming. *Proceedings of the National Academy of*
53 *Sciences of the United States of America*, 103(1):129–134, 2006.
54
55
56 [37] S. S. Óhéigeartaigh, D. Armisén, K. P. Byrne, and K. H. Wolfe. SearchDOGS bacteria, software that
57 provides automated identification of potentially missed genes in annotated bacterial genomes. *Journal*
58 *of bacteriology*, 196(11):2030–2042, 2014.
59
60
61
62
63
64
65

- 1
2
3
4 [38] K. Osatomi and H. Sumiyoshi. Complete nucleotide sequence of dengue type 3 virus genome RNA.
5 *Virology*, 176(2):643 – 647, 1990.
6
7
8 [39] R. Percudani, D. Carnevali, and V. Puggioni. Ureidoglycolate hydrolase, amidohydrolase, lyase: how
9 errors in biological databases are incorporated in scientific papers and vice versa. *Database*, 2013:bat071,
10 2013.
11
12
13 [40] F. Pfeiffer and D. Oesterhelt. A manual curation strategy to improve genome annotation: Application
14 to a set of haloarchael genomes. *Life*, 5(2):1427–1444, 2015.
15
16
17 [41] M. S. Poptsova and J. P. Gogarten. Using comparative genome analysis to identify problems in annotated
18 microbial genomes. *Microbiology*, 156(7):1909–1917, 2010.
19
20
21 [42] S. Poux, M. Magrane, C. N. Arighi, A. Bridge, C. O’Donovan, K. Laiho, U. Consortium, et al. Ex-
22 pert curation in uniprotkb: a case study on dealing with conflicting and erroneous data. *Database*,
23 2014:bau016, 2014.
24
25
26 [43] V. J. Promponas, I. Iliopoulos, and C. A. Ouzounis. Annotation inconsistencies beyond sequence
27 similarity-based function prediction – phylogeny and genome structure. *Standards in Genomic Sciences*,
28 10(1):108, 2015.
29
30
31
32 [44] P. Radivojac, W. T. Clark, T. R. Oron, A. M. Schnoes, T. Wittkop, A. Sokolov, K. Graim, C. Funk,
33 K. Verspoor, A. Ben-Hur, et al. A large-scale evaluation of computational protein function prediction.
34 *Nature methods*, 10(3):221–227, 2013.
35
36
37 [45] S. E. Robertson, S. Walker, S. Jones, M. Hancock-Beaulieu, and M. Gatford. Okapi at trec-2. In *TREC*,
38 pages 21–34, 1993.
39
40
41 [46] A. Rudniy, M. Song, and J. Geller. Detecting duplicate biological entities using shortest path edit
42 distance. *Int. J. Data Min. Bioinformatics*, 4(4):395–410, July 2010.
43
44
45 [47] E. Sayers. E-utilities quick start. entrez programming utilities help. Technical report, 2010.
46
47
48 [48] D. Schatz, A. Shemi, S. Rosenwasser, H. Sabanay, S. G. Wolf, S. Ben-Dor, and A. Vardi. Corrigendum.
49 *New Phytologist*, 206(2):881–881, 2015.
50
51
52 [49] A. M. Schnoes, S. D. Brown, I. Dodevski, and P. C. Babbitt. Annotation error in public databases:
53 Misannotation of molecular function in enzyme superfamilies. *PLoS Comput Biol*, 5(12):1–13, 12 2009.
54
55
56 [50] G. A. Seluja, A. Farmer, M. McLeod, C. Harger, and P. A. Schad. Establishing a method of vector
57 contamination identification in database sequences. *Bioinformatics*, 15(2):106–110, 1999.
58
59
60 [51] R. J. Siezen and S. A. Van Hijum. Genome (re-) annotation and open-source annotation pipelines.
61 *Microbial biotechnology*, 3(4):362–369, 2010.
62
63
64
65

- 1
2
3
4 [52] M. Song and A. Rudniy. Detecting duplicate biological entities using markov random field-based edit
5 distance. In *Bioinformatics and Biomedicine, 2008. BIBM '08. IEEE International Conference on*,
6 pages 457–460, Nov 2008.
7
8
9
10 [53] K. Srinivasan, P. Gopalakrishnakone, P. Tan, K. Chew, B. Cheng, R. Kini, J. Koh, S. Seah, and
11 V. Brusica. Scorpion, a molecular database of scorpion toxins. *Toxicon*, 40(1):23 – 31, 2002.
12
13 [54] *The Gene Ontology Consortium*. Gene ontology: tool for the unification of biology. *Nat Genet*, 25:25
14 – 29, 2000.
15
16 [55] A. Tritt, J. A. Eisen, M. T. Facciotti, and A. E. Darling. An integrated pipeline for de novo assembly
17 of microbial genomes. *PloS one*, 7(9):e42304, 2012.
18
19 [56] Q. Wu, Y. Ye, M. K. Ng, S.-S. Ho, and R. Shi. Collective prediction of protein functions from protein-
20 protein interaction networks. *BMC bioinformatics*, 15(2):1, 2014.
21
22 [57] R. Zallot, K. J. Harrison, B. Kolaczowski, and V. de Crécy-Lagard. Functional annotations of paralogs:
23 A blessing and a curse. *Life*, 6(3):39, 2016.
24
25 [58] C. Zhai and J. Lafferty. A study of smoothing methods for language models applied to ad hoc information
26 retrieval. In *Proceedings of the 24th Annual International ACM SIGIR Conference on Research and*
27 *Development in Information Retrieval, SIGIR '01*, pages 334–342, New York, NY, USA, 2001. ACM.
28
29 [59] Y. Zhao, F. Scholer, and Y. Tsegay. Effective pre-retrieval query performance prediction using similarity
30 and variability evidence. In *30th European Conference on IR Research, ECIR '08*, pages 52–64, Berlin,
31 Heidelberg, 2008. Springer Berlin Heidelberg.
32
33 [60] E. Zorita, P. Cusco, and G. J. Filion. Starcode: sequence clustering based on all-pairs search. *Bioin-*
34 *formatics*, 31(12):1913–1919, 2015.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65