

This is an authors' manuscript version of a paper accepted for proceedings of *TPDL-2018, Porto, Portugal, Sept 10-13*. The final authenticated publication is available online at <https://doi.org/>will be added as soon as available

Content-Based Quality Estimation for Automatic Subject Indexing of Short Texts under Precision and Recall Constraints

Martin Toepfer¹ and Christin Seifert²

¹ ZBW – Leibniz Information Centre for Economics, Kiel, Germany
m.toepfer@zbw.eu

² University of Twente, Enschede, The Netherlands
c.seifert@utwente.nl

Abstract. Semantic annotations have to satisfy quality constraints to be useful for digital libraries, which is particularly challenging on large and diverse datasets. Confidence scores of multi-label classification methods typically refer only to the relevance of particular subjects, disregarding indicators of insufficient content representation at the document-level. Therefore, we propose a novel approach that detects documents rather than concepts where quality criteria are met. Our approach uses a deep, multi-layered regression architecture, which comprises a variety of content-based indicators. We evaluated multiple configurations using text collections from law and economics, where the available content is restricted to very short texts. Notably, we demonstrate that the proposed quality estimation technique can determine subsets of the previously unseen data where considerable gains in document-level recall can be achieved, while upholding precision at the same time. Hence, the approach effectively performs a filtering that ensures high data quality standards in operative information retrieval systems.

Keywords: Quality Estimation, Automatic Subject Indexing, Document-Level Constraints, Multi-Label Classification, Meta-Learning, Short-Text

1 Introduction

Semantic annotations from automatic subject indexing can improve information retrieval (IR) by query expansion, however, classification performance is a critical factor to gain the benefits [1]. The relevance of multi-label text classification engendered research in several disciplines. Although considerable progress has been made over the last decades [2,3,4,5,6,7,8,9], several challenges remain. Just to give an example, $\text{precision@5} = 52\%$ [9] has recently been reported for a dataset in the legal domain (EURLEX [5]), which means that on average per document only half of the five top-ranked subjects matched human annotations. Institutional quality requirements, like for instance at libraries, often put severe constraints on precision [10] as well as recall. It is therefore not sufficient to just apply state-of-the-art algorithms with respect to averaged f_1 scores, but further

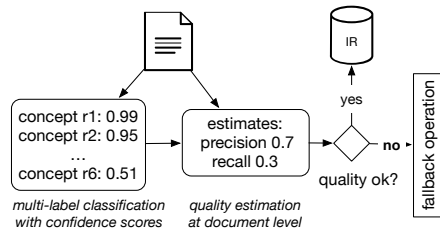


Fig. 1: Schematic overview of the main application context. Document-level quality estimation enables filtering of automatic subject indexing results.

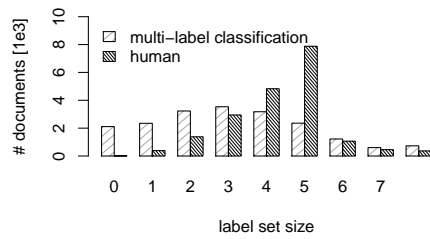


Fig. 2: Illustration of low document-level recall by comparing distributions of label set size (human vs. multi-label classification) [Dataset: EURLEX].

necessary to automatically separate the wheat from the chaff. On large datasets, human quality assessment can be too expensive. While specific confidence estimation approaches have been proposed for domains like information extraction [11], automatic subject indexing and multi-label text classification miss essential research in this direction. We aim to fill this gap since quality estimation is becoming increasingly important for digital libraries to integrate autonomous processes into operative IR systems. As depicted in Fig. 1, document-level quality estimates allow to implement filters at the interface to databases which are used for IR.

Most automatic subject indexing methods provide a score for each concept [2,12,6,7], hence allowing to exclude individual predictions that might be incorrect. Such precision-oriented filtering removes single label assignments from documents, leading to lower document-level recall, as exemplified in Fig. 2. As a direct consequence, it becomes difficult to assess document-level quality. As can be seen, the plain number of assigned concepts to a document is not a satisfying indicator, since human indexers³ use a wide range of label set sizes. Interestingly, Sect. 3 will point out that uncertainty in recall is an inherent and inevitable phenomenon of multi-label text classification when only a few preconditions are met. Thus we conclude that concept-level confidence scores must be complemented with document-level estimates, as investigated in this paper.

In summary, the contributions of this work are the following:

- We provide a conceptual analysis of confidence and quality estimation for automatic subject indexing.
- We propose a quality estimation approach, termed *Qualle*, which combines multiple content-based features in a multi-layered regression architecture.
- We show the impact of different feature groups and the effectiveness of *Qualle* for quality estimation and filtering in an empirical study.

³ For brevity, the remainder of this paper will simply use the terms indexing, indexer, ... to refer to subject indexing, subject indexer, ..., respectively.

The remainder of the paper starts with a discussion of related work before the central section introduces the quality estimation approach (Sect. 3) followed by experimental results (Sect. 4).

2 Related Work

Confidence scores are an integral part of many machine learning (ML) approaches for multi-label text classification [2,13]. For instance, rule-learning typically computes a confidence score for each rule, dividing the number of times the rule correctly infers a class label by the number of times the rule matches in total. Naive-Bayes approaches use Bayes' Rule to derive conditional probabilities. Flexible techniques have been developed to perform probability calibration [14]. Thus, systems using multi-label classification (MLC) machine learning methods for subject indexing often provide confidence scores for each subject heading. [2,13]. Medelyan and Witten [12] used decision trees to compute confidence scores for dictionary matches. Huang et al. [6] similarly applied a learning-to-rank approach on MeSH term recommendation based on candidates from k-nearest-neighbors. In general, binary relevance (BR) approaches also provide probabilities for each concept, for instance by application of probability calibration techniques (e.g. [7]). Tang et al. [15] proposed a BR system which additionally creates a distinct model to determine the number of relevant concepts per document. In summary, the scores provided by the above mentioned systems are limited to concept-level confidence, that is, referring to distinct subjects.

In the context of classifier combination, Bennett et al. [3] proposed *reliability-indicator variables* for model selection. They identified four types of indicator variables and showed their utility. In contrast to their work, we focus on different objectives. We apply such features (reliability indicators) for quality estimation, which in particular comprises estimation of recall. By contrast, *precision-constrained* situations have recently been studied by Bennett et al. [10]. *Confidence in predictions and classifiers* has recently gained attention in the context of *transparent machine learning* (e.g. [16]). Contrary to transparent machine learning, quality estimation does not aim to improve interpretability, and it thus may be realized by black box ML models. Nevertheless, quality estimates may be relevant for humans to gain trust in ML.

Confidence estimation has been studied in different *application domains*, and it has been noted that different levels of confidence scores are relevant. For instance, Culotta and McCallum [11] distinguished between field confidence and record confidence (entire record is labeled correctly) in information extraction. They compared different scoring methods and also trained a classifier that discriminates correct and incorrect instances for fields and records, respectively.

3 Quality Estimation

Our approach to quality estimation (Sect. 3.2) stems from an analysis of common practice, as described in the following.

3.1 Analysis

In the past, quality of automatic subject indexing has been assessed in different ways that have individual drawbacks. Traditionally, library and information scientists regarded indexing quality, effectiveness, and consistency [17]. Quality assessment that requires human judgements is, however, costly, which can be a severe issue on large and diverse datasets. For this reason, evaluations of automatic subject indexing often just rely on consistency with singly annotated human indexing, yielding metrics which are known as precision and recall. As described in Sect. 2, common indexing approaches provide confidence scores for each class, denoting posterior probabilities $p(y_j = 1|\mathcal{D})$, where y_j refers to a single concept of the controlled vocabulary, thus they are referred to as *concept-level confidence* in this work. Statistical associative approaches derive confidence scores based on dependencies between terms and class labels from examples. As a consequence, the performance of these methods largely depends on the availability of appropriate training examples and the stability of term and concept distributions, whereas lexical methods require vocabularies that exhaustively cover the domain. When concept drift occurs, that is, if observed terms and the set of relevant concepts differ between training data and new data, both types of indexing approaches considerably decrease in performance [8]. Interestingly, since these algorithms merely learn to assign recognized subjects of the controlled vocabulary, they will silently miss to assign relevant subjects not covered by the controlled vocabulary, and moreover they are unable to recognize and represent the loss in document-level content representation. It is further plausible that these issues are more pronounced when only titles of documents are processed, since for title-based indexing the complete subject content is compressed into only a few words which makes understanding of each single word more crucial compared to processing full texts. As the evolution of terms and concepts is an inherent property of language (cf. e.g. [18]), accurate recognition of insufficient exhaustivity is essential in the long term. It must be assumed that uncertainty in recall is an inherent and inevitable phenomenon of automatic indexing and multi-label text classification in general. For these reasons, in order to guarantee quality, indexing systems must gain knowledge relating to classifier reliability based on additional representations (cf. [3]), exploiting information such as out-of-vocabulary term occurrences and document length, just to give an example. Therefore, instead of concept-level confidence, we propose to address document-level quality directly.

3.2 Quelle: Content-Based Quality Estimation

Multi-label classification methods can be tuned by regularization and configuration of thresholds to satisfy constraints on precision. Hence, the main challenge for our approach on quality estimation, *Quelle*, is to estimate document-level recall. As a solution, we propose the architecture which is exemplified in Fig. 3. The input layer shows a fictitious title of a document to be indexed, which is then represented by multiple features. The content is processed by a multi-label

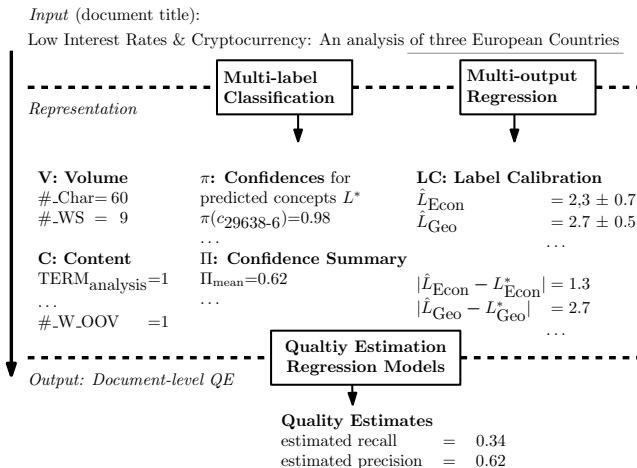


Fig. 3: Multi-layered regression architecture for quality estimation (example).

text classification method, producing a set of concepts and corresponding confidence values, e.g., $\pi(c_{29638-6}) = 0.98$.⁴ Moreover, a multi-output regression module offers expectations regarding the proper number of concepts for a document (Label Calibration), if possible considering distinct semantic categories, for instance geographic names (\hat{L}_{Geo}) or economics subject terms (\hat{L}_{Econ}), commodities, and much more. In the given example, the phrase “three European countries” clearly points out that it would be reasonable to assume three geographic names when access to the full text is possible, however, without particularly specifying which ones to choose. The input is not precise enough. Drawing connections between the predicted concept set L^* and the estimated numbers of concepts \hat{L} can indicate such recall issues. For instance, $|\hat{L}_{Geo} - L_{Geo}^*| = 2.7$ indicates that the proposed index terms miss more than two geographic names. Such reasoning is not covered by ordinary statistical text categorization methods. In addition, basic reliability indicators are included as features, such as content length (#.Char), individual term indicators (e.g.: TERM_{analysis}), the number of out-of-vocabulary terms (#.W_OOV), or different types of aggregations (Π) of the confidence scores of the assigned concepts. Finally, quality aspects are estimated using regression models based on the aforementioned features.

Development of the feature groups (Fig. 3: V, C, Π , LC) was driven by conceptual considerations. In particular, we wanted the features to represent: *imprecise input* (e.g., “three European countries”: inherent ambiguity), *lack of input information* (e.g., title with fewer than 4 words: information is scarce), as well as *lack of knowledge* (e.g., “On Expected Effects of Brexit on European law”: information is present but can not be interpreted, if the term “Brexit” has not been observed before).

⁴ The concept identifier 29638-6 refers to the concept “Low-interest-rate policy”.

In general, the architecture of Qualle is a framework which, for example, allows to apply arbitrary regression methods for quality estimation. Since the number of completely correct records in automatic subject indexing is extremely low, we do not consider re-ranking by MaxEnt, which has been investigated for record-based confidence estimation in information extraction [11]. In this paper, we focus our analysis of the deep quality estimation regression approach on document-level recall. In addition, basic indicators have been considered for document-level precision estimation, that is, the mean (II_{mean}), product, median and minimum of the confidence values of the assigned concepts.

4 Experiments

The experiments are centered around the following questions:

- Q1: Do predictions of recall and actual recall correlate with each other?
- Q2: How accurate are the recall estimates?
- Q3: Which of the feature groups contribute most to recall prediction?
- Q4: What are the effects of filtering based on recall estimates on coverage and document-level precision?

which are relevant in practice. Ranking documents by document-level recall (Q1) allows to separate high-recall documents from low-recall documents. Accurate estimates (Q2) allow to control filtering with meaningful constraints. Applicability of the filtering approach would, however, be prevented if either document-level precision was decreased considerably or if the number of documents passing the filter was too low (Q4).

4.1 Setup

We evaluate the approach in two domains. We first perform a basic experiment on legal texts, addressing questions Q1 and Q2. Subsequently, we go into details regarding economics literature, treating questions Q1-Q4. The second experiments use more, and more elaborated features that notably exploit the categorisation of concepts into a hierarchy, as described in Section 3.2.

The adequacy of quality estimation is measured in two ways. Since perfect quality estimates follow their corresponding actual counterparts linearly, we consider the Pearson product-moment *correlation* coefficient ρ for (Q1)⁵. A strong correlation between predicted and true quality allows to order documents correctly, that is, corresponding to the true performance. ρ has been used in related studies [11]. For measuring the exactness of estimated recall values (Q2), we consider the *mean squared error* (MSE). To gain knowledge about the utility of the feature groups (Q3), we perform a systematic analysis of different configurations. Feature groups are removed separately from the complete set of features (ablation study), and measurements are also collected for each feature group alone (isolation study). Question Q4 was addressed by evaluating different thresholds

⁵ If only ranking is relevant, rank-based correlation coefficients should be considered.

on estimated recall and measuring average true precision and recall over the corresponding selected documents. In addition, the *coverage* = $\frac{|D_{\text{selected}}|}{N}$ was measured, with N being the total number of documents and D_{selected} the selected subset of the whole data set. We also report the relative recall *gain* (RG) on these subsets. The accuracy of initial multi-label classification is reported briefly for comparability, using metrics as described in Sect. 3.

Regarding law, we employ the EURLEX [5] dataset to address Q1 and Q2. It comprises 19,314 documents, each having 5.31 EUROVOC⁶ subject terms on average. For further details on the data set, refer to [5] and the website of the dataset⁷. Please note that the experiments in this paper only use the titles rather than the full text of the documents and that different train/test splits were used. Regarding economics, we use three datasets, which comprise roughly 20,000 (T20k), 60,000 (T60k), and 400,000 documents (T400k), respectively. Each document is associated with several descriptors, for instance 5.89 on average for T400k, from the STW Thesaurus for Economics (STW)⁸. Both, the STW and EUROVOC, comprise thousands of concepts, yielding challenging multi-label classification tasks.

For each data set, we perform cross validation with 5-folds. And for each of those 5 runs, we apply nested cross validation runs, likewise with 5 folds used for parameter optimization and learning of quality estimation. That is, each training set is subdivided into *dev-train* and *dev-test* splits. For validation, a new model is trained from random samples of the same size as one of the dev-train splits. As a consequence, the training and prediction of the classifier for label prediction as well as for the regressor for label calibration are carried out $5 \cdot 6 = 30$ times for each collection. Quality estimation is evaluated on the corresponding *eval-test* data folds.

For *multi-label text classification*, we chose binary relevance logistic regression (BRLR) optimized with stochastic gradient descent (cf. [7,10]).

Regarding *reliability indicator variables*, the EURLEX study relies on just two features: the estimated number of concepts for the document, and the difference to the actually predicted number of concepts for the document by BRLR. For the detailed study on economics documents, all feature groups were employed (Sect. 3.2). Label calibration has been realized with tree-based methods (EURLEX: ExtraTreesRegressor [19], Economics: GradientBoostingRegressor [20]). Only the total number of concepts per document is considered for EURLEX. The economics experiments compute label calibration estimates for the seven top categories of the corresponding thesaurus. For EURLEX and economics, *#_Char*, *#_WS* and *TERM_i* have been used as features for label calibration.

Several regression methods implemented in scikit-learn [21] were considered for quality estimation. For the EURLEX experiments, rather basic models like LinearRegression and DecisionTreeRegression are tested, as well as ensemble machine learning methods, namely, ExtraTrees [19], GradientBoosting [20], and

⁶ <http://eurovoc.europa.eu/>, accessed: 31.12.2017

⁷ <http://www.ke.tu-darmstadt.de/resources/eurllex>, accessed 31.12.2017

⁸ <http://zbw.eu/stw/version/latest/about.en.html>, accessed: 09.01.2018

AdaBoostRegressor [22]. different *regression methods* were applied for recall estimation. Regarding the more detailed experiments on economics research literature, only the two regression models that performed best on EURLEX were investigated. Extensive grid searches over various parameters of the models are left for future work.

4.2 Results

EURLEX From the different regression models, LinearRegression produced the lowest correlation coefficient ($\rho = .214 \pm .026$) between predicted recall and true recall. AdaBoostRegressor reached the highest correlation coefficient ($\rho = .590 \pm .013$) and the lowest mean squared error ($\text{MSE} = 0.067 \pm 0.002$). Only AdaBoostRegressor and GradientBoostingRegressor achieved correlation coefficients greater than .500. Although being worse than the AdaBoostRegressor on average, the results for the ExtraTreesRegressor were more balanced.

Economics Comparing the two selected regression methods, we found that the best configurations of GradientBoosting dominated the best configurations of AdaBoost on all datasets and with respect to both metrics (ρ , MSE). Thus, Adaboost has been excluded from further analysis.

Table 1 offers the numbers for ablation and isolation of feature groups. For each collection, the complete set of features (first row corresponding to each collection) is always among the top configurations, where differences are not greater than the sum of their standard deviations. For all collections, the largest decrease in performance is recognized when the group of features related to label calibration is removed. In accordance, this feature group yields the strongest individual results. On T20k, its performance is close to the complete set of features. For the collections with more data, the difference is more clear. Volume features, including length of the document, was found to be the lowest ranking group and has little impact when removed from the complete set of features. In nearly all cases of configurations, more data yields higher correlation coefficients, however, not necessarily lower mean squared error. In the following, we focus on reporting results regarding T400k. Figures for T20k and T60k were similar.

Figure 4a) depicts recall estimation results for T400k. The plot illustrates the degree of linear relation and also reveals the distributions of estimated and true recall values. Most of the documents have a true recall that is less than 60%. Regarding the scoring functions for document-level precision, the product of concept-level confidence scores exhibited the highest correlations for T20k and T60k, however, still staying below .500. On T400k, all scoring functions were very close to each other, and their correlation coefficients were above .500. Figure 4b) depicts results for the product of concept confidence values.

Finally, Fig. 5 visualizes how different thresholds on estimated recall affect properties of the resulting document selections. The plot therefore shows coverage, as well as mean document-level true recall and true precision. When constraining estimated recall to be at least 30%, a gain $\text{RG}=44\%$ of true recall in

Table 1: Feature analysis for economics with GradientBoosting. \checkmark : presence of feature group. Δ : Difference in relation to complete set of features. \dagger : Absolute difference to condition with all features is greater than the sum of their *sd*.

Configuration	V	C	LC	Π	$\rho \pm$ std	Δ_ρ	MSE \pm std	Δ_{MSE}
T20k	\checkmark	\checkmark	\checkmark	\checkmark	0.597 ± 0.014	-0.0%	0.039 ± 0.001	-0.0%
T20k		\checkmark	\checkmark	\checkmark	0.596 ± 0.014	-0.2%	0.040 ± 0.001	0.2%
T20k	\checkmark		\checkmark	\checkmark	0.595 ± 0.015	-0.3%	0.039 ± 0.001	-0.6%
T20k	\checkmark	\checkmark	\checkmark		0.583 ± 0.015	-2.3%	0.040 ± 0.001	1.8%
T20k	\checkmark	\checkmark		\checkmark	0.384 ± 0.005	-35.6% \dagger	0.050 ± 0.001	26.5% \dagger
T20k			\checkmark		0.569 ± 0.014	-4.7% \dagger	0.041 ± 0.001	2.6%
T20k				\checkmark	0.362 ± 0.007	-39.3% \dagger	0.051 ± 0.001	28.0% \dagger
T20k		\checkmark			0.196 ± 0.013	-67.1% \dagger	0.056 ± 0.001	41.1% \dagger
T20k	\checkmark				0.128 ± 0.008	-78.6% \dagger	0.056 ± 0.001	43.0% \dagger
T60k	\checkmark	\checkmark	\checkmark	\checkmark	0.617 ± 0.011	-0.0%	0.043 ± 0.000	-0.0%
T60k		\checkmark	\checkmark	\checkmark	0.615 ± 0.010	-0.3%	0.044 ± 0.000	0.3%
T60k	\checkmark		\checkmark	\checkmark	0.602 ± 0.009	-2.5%	0.044 ± 0.001	1.8%
T60k	\checkmark	\checkmark	\checkmark		0.600 ± 0.010	-2.8%	0.044 ± 0.000	2.4% \dagger
T60k	\checkmark	\checkmark		\checkmark	0.420 ± 0.009	-31.9% \dagger	0.055 ± 0.001	26.1% \dagger
T60k			\checkmark		0.574 ± 0.005	-6.9% \dagger	0.046 ± 0.001	5.4% \dagger
T60k				\checkmark	0.391 ± 0.011	-36.6% \dagger	0.056 ± 0.001	28.7% \dagger
T60k		\checkmark			0.216 ± 0.017	-64.9% \dagger	0.062 ± 0.001	43.9% \dagger
T60k	\checkmark				0.069 ± 0.009	-88.8% \dagger	0.064 ± 0.001	48.2% \dagger
T400k	\checkmark	\checkmark	\checkmark	\checkmark	0.648 ± 0.002	-0.0%	0.042 ± 0.000	-0.0%
T400k	\checkmark	\checkmark	\checkmark		0.649 ± 0.001	0.1%	0.042 ± 0.000	-0.1%
T400k		\checkmark	\checkmark	\checkmark	0.648 ± 0.001	0.0%	0.042 ± 0.000	0.2%
T400k	\checkmark		\checkmark	\checkmark	0.644 ± 0.002	-0.6% \dagger	0.042 ± 0.000	0.8% \dagger
T400k	\checkmark	\checkmark		\checkmark	0.528 ± 0.002	-18.5% \dagger	0.050 ± 0.000	19.5% \dagger
T400k			\checkmark		0.640 ± 0.001	-1.3% \dagger	0.043 ± 0.000	1.1% \dagger
T400k				\checkmark	0.511 ± 0.002	-21.2% \dagger	0.051 ± 0.000	21.6% \dagger
T400k		\checkmark			0.225 ± 0.003	-65.3% \dagger	0.064 ± 0.000	51.6% \dagger
T400k	\checkmark				0.122 ± 0.002	-81.1% \dagger	0.065 ± 0.000	55.3% \dagger

relation to the measure on the complete collection could be achieved on T400k. The most relevant message that can be drawn from Fig. 5 is that the precision on the selected subsets remained the same or even increased, when putting harder constraints on estimated recall.

Multi-label Classification The performance of the BRLR approach was not in the focus of the study, yet BRLR turned out to be a reasonable choice. For instance, it reached sample-based average $f_1 = 0.361$, precision = 0.528, recall = 0.327 on T20k, and $f_1 = 49.1\%$ (micro avg.) on EURLEX. These figures broadly correspond to related studies.

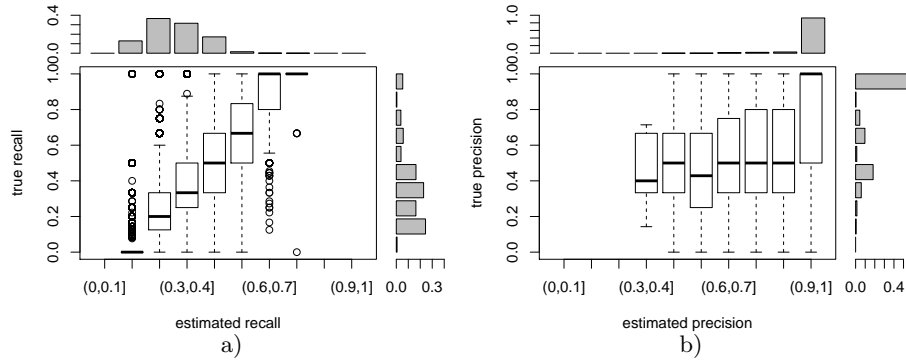


Fig. 4: Quality estimates and true values (Economics: T400k). a) Recall estimation by Qualle, b) Precision score by product of concept confidence values. Marginal distributions (bin count/total count) are shown on the top and on the right, respectively.

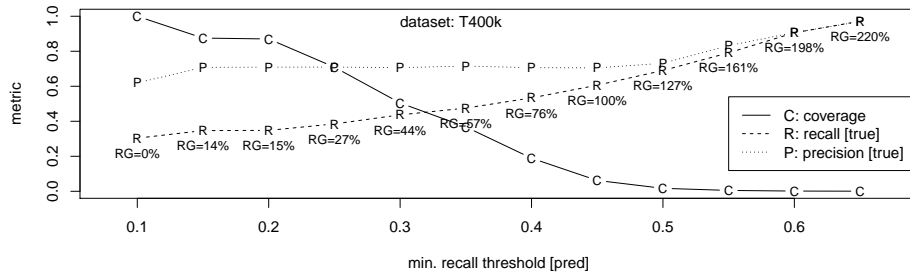


Fig. 5: Quality vs. coverage (collection: T400k): coverage, mean document-level true recall and true precision for different predicted recall thresholds. RG: relative gain in document-level recall on selected subset compared to the full dataset.

4.3 Discussion

The basic set of features used on EURLEX reached respectable correlations ($\rho > .500$) between predicted and true recall (Q1) only for the sophisticated machine learning methods AdaBoost and GradientBoosting. Differences in the balance of predictions should be considered for applications, just like the notable amount of variance that remains around the predicted values (Q2). In summary, the outcome of the EURLEX study suggested that recall estimates that are useful for filtering are feasible, motivating investigation of more complex configurations.

Looking at the outcome of the experiments on the economics datasets, especially Fig. 5, our results show that the proposed quality estimation approach can be successfully applied to identify subsets of document collections where soft constraints on precision as well as recall are met (Q4). Finally, it remains a decision depending on the application context to make trade-offs according

to multi-criteria objectives, which notably comprise coverage. Regarding recall, ranking and accuracy of predictions are sufficient enough (Q1, Q2). Interestingly, precision was not affected negatively (cf. Fig. 5). Based on Table 1, applications should consider the full set of features, which belongs to the top performing configurations in all cases and outperformed individual feature groups. Label calibration information is found to be a strong individual predictor. It is the most relevant reliability indicator (Q3) compared to the volume, content, and concept-confidence related feature groups. The mean squared errors of predictions indicate that a considerable amount of vagueness remains (Q2). Possibly, it may be caused by the errors in concept assignments, which influence the label calibration related features.

The experimental results highlight the inevitable difficulties (cf. Sect. 3.1) in multi-label text classification, namely, suffering from low document-level recall when the model misses knowledge (either dictionary entries or training examples), or when the observed input is inherently ambiguous. Quality estimation enables to handle such issues by controlling, that is, making trade-offs between quality and coverage. Since the proposed approach is not bound to specific MLC or regression methods, further progress in this regard can be integrated and is assumed to improve collection coverage. Another direction for future work is to consider alternative quality metrics that take semantic relations into account (see e.g., [4,23]).

5 Conclusion

In order to assure data quality in operative information retrieval systems with large and diverse datasets, we investigated an important yet less addressed research topic, namely quality estimation of automatic subject indexing with a focus on the document level. Our experimental results on two domains spanning over collections of different sizes show that the proposed multi-layer architecture is effective and thus enables to control quality in settings where high standards have to be met. The approach allows to define different thresholds, which resulted in considerable gains of document-level recall, while upholding precision at the same time. Label calibration was the most relevant reliability indicator.

References

1. Trieschnigg, D., Pezik, P., Lee, V., de Jong, F., Kraaij, W., Rebholz-Schuhmann, D.: MeSH up: effective MeSH text classification for improved document retrieval. *Bioinformatics* **25**(11) (2009) 1412–1418
2. Sebastiani, F.: Machine learning in automated text categorization. *ACM Computing Surveys* **34**(1) (2002) 1–47
3. Bennett, P.N., Dumais, S.T., Horvitz, E.: Probabilistic combination of text classifiers using reliability indicators: models and results. In Jävelin, K., Beaulieu, M., Baeza-Yates, R.A., Myaeng, S., eds.: *Proc. SIGIR 2002*, ACM (2002) 207–214
4. Medelyan, O., Witten, I.H.: Measuring inter-indexer consistency using a thesaurus. In: *Proc. JCDL 2006*, ACM (2006) 274–275

5. Loza Mencía, E., Fürnkranz, J.: Efficient multilabel classification algorithms for large-scale problems in the legal domain. In: *Semantic Processing of Legal Texts – Where the Language of Law Meets the Law of Language*. Volume 6036 of LNAI. 1st edn. Springer (2010) 192–215
6. Huang, M., Névéol, A., Lu, Z.: Recommending MeSH terms for annotating biomedical articles. *JAMIA* **18**(5) (2011) 660–667
7. Wilbur, W.J., Kim, W.: Stochastic gradient descent and the prediction of MeSH for PubMed records. *AMIA Annu Symp Proc* **2014** (2014) 1198–1207
8. Toepfer, M., Seifert, C.: Descriptor-invariant fusion architectures for automatic subject indexing. In: *Proc. JCDL 2017*, IEEE Computer Society (2017) 31–40
9. Liu, J., Chang, W., Wu, Y., Yang, Y.: Deep learning for extreme multi-label text classification. In Kando, N., Sakai, T., Joho, H., Li, H., de Vries, A.P., White, R.W., eds.: *Proc. SIGIR 2017*, ACM (2017) 115–124
10. Bennett, P.N., Chickering, D.M., Meek, C., Zhu, X.: Algorithms for active classifier selection: Maximizing recall with precision constraints. In de Rijke, M., Shokouhi, M., Tomkins, A., Zhang, M., eds.: *Proc. WSDM 2017*, ACM (2017) 711–719
11. Culotta, A., McCallum, A.: Confidence estimation for information extraction. In: *Proc. HLT-NAACL 2004: Short Papers*, ACL (2004) 109–112
12. Medelyan, O., Witten, I.H.: Domain-independent automatic keyphrase indexing with small training sets. *JASIST* **59**(7) (2008) 1026–1040
13. Gibaja, E., Ventura, S.: A tutorial on multilabel learning. *ACM Comput. Surv.* **47**(3) (2015) 52:1–52:38
14. Zadrozny, B., Elkan, C.: Transforming classifier scores into accurate multiclass probability estimates. In: *Proc. SIGKDD 2002*, ACM (2002) 694–699
15. Tang, L., Rajan, S., Narayanan, V.K.: Large scale multi-label classification via metalabeler. In: *Proc. WWW 2009*, ACM (2009) 211–220
16. Ribeiro, M.T., Singh, S., Guestrin, C.: Why should I trust you?: Explaining the predictions of any classifier. In: *Proc. SIGKDD 2016*, ACM (2016) 1135–1144
17. Rolling, L.N.: Indexing consistency, quality and efficiency. *Information Processing & Management* **17**(2) (1981) 69–76
18. Tahmasebi, N., Risse, T.: On the uses of word sense change for research in the digital humanities. In Kamps, J., Tsakonas, G., Manolopoulos, Y., Iliadis, L.S., Karydis, I., eds.: *Proc. TPD L 2017*. Volume 10450 of LNCS., Springer (2017) 246–257
19. Geurts, P., Ernst, D., Wehenkel, L.: Extremely randomized trees. *Machine Learning* **63**(1) (Apr 2006) 3–42
20. Friedman, J.H.: Stochastic gradient boosting. *Computational Statistics & Data Analysis* **38**(4) (2002) 367–378
21. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12** (2011) 2825–2830
22. Drucker, H.: Improving regressors using boosting techniques. In Fisher, D.H., ed.: *Proc. ICML 1997*, Morgan Kaufmann (1997) 107–115
23. Neveol, A., Zeng, K., Bodenreider, O.: Besides precision & recall: exploring alternative approaches to evaluating an automatic indexing tool for medline. *AMIA Annu Symp Proc* (2006) 589–593