

Automatic Extension of Corpora from the Intelligent Ensembling of eHealth Knowledge Discovery Systems Outputs

Juan Pablo Consuegra-Ayala^{a,*}, Yoan Gutiérrez^{b,c}, Alejandro Piad-Morffis^a, Yudivian Almeida-Cruz^a, Manuel Palomar^{b,c}

^a*School of Math and Computer Science, University of Habana, La Habana, Cuba, 10200.*

^b*University Institute for Computing Research (IUII), University of Alicante, Alicante, Spain, 03690.*

^c*Department of Language and Computing Systems, University of Alicante, Alicante, Spain, 03690.*

Abstract

Corpora are one of the most valuable resources at present for building machine learning systems. However, building new corpora is an expensive task, which makes the automatic extension of corpora a highly attractive task to develop. Hence, finding new strategies that reduce the cost and effort involved in this task, while at the same time guaranteeing quality, remains an open and important challenge for the research community. In this paper, we present a set of ensembling strategies oriented toward entity and relation extraction tasks. The main goal is to combine several automatically annotated versions of corpora to produce a single version with improved quality. An ensembler is built by exploring a configuration space in search of the combination that maximizes the fitness of the ensembled collection according to a reference collection. The eHealth-KD 2019 challenge was chosen for the case study. The submitted systems' outputs were ensembled, resulting in the construction of

*Corresponding author.

Email addresses: jpconsuegra@matcom.uh.cu (Juan Pablo Consuegra-Ayala), ygutierrez@dlsi.ua.es (Yoan Gutiérrez), apiad@matcom.uh.cu (Alejandro Piad-Morffis), yudy@matcom.uh.cu (Yudivian Almeida-Cruz), mpalomar@dlsi.ua.es (Manuel Palomar)

an automatically annotated collection of 8 000 sentences. We show that using this collection as additional training input for a baseline algorithm has a positive impact on its performance. Additionally, the ensembling pipeline was used as a participant system in the 2020 edition of the challenge. The ensembled run achieved a slightly better performance than the individual runs.

Keywords:

Ensemble Methods, Annotated Corpora, Information Extraction, Entity Recognition, Relation Extraction, Natural Language Processing

1. Introduction

In recent years, there has been an increasing interest in applying machine learning and Natural Language Processing (NLP) techniques to solve several types of complex tasks. Named entity recognition (NER), relation extraction, machine translation, and sentiment analysis are just a few well-known examples of these tasks [1, 2, 3, 4]. These and other NLP tasks have a considerable impact on society nowadays and machine learning techniques have been able to achieve good performance [5, 6]. Although systems based on handcrafted rules might be able to perform fairly well in some of these tasks [7], the use of machine learning models eliminates the need for expert knowledge when designing new rules, and this reduces the cost of solving new problems.

One of the major issues involved in applying machine learning technologies is that they require large collections of corpora for building learning models. In some cases, the corpora are a natural component of the problem and therefore can be easily obtained by collecting known examples. That is the case with more classical approaches, such as stock market predictions based on features computed from automatically collected data [8]. However, in other cases, the corpora must be

18 handcrafted to establish a training base for solving certain subtasks. This frequently
19 occurs with information extraction problems, where specific annotations schemes are
20 defined and then extracted from plain text [9].

21 These corpora are considered highly valuable resources nowadays, not only because
22 they are highly demanded, but they are also difficult and expensive to build. Hence,
23 techniques have been developed for reducing the complexity of building corpora.
24 Some technologies assist the corpus annotation directly, making it easier for humans
25 to manually build the corpus. Such systems have applied active learning techniques to
26 train a tagging model [10]. The human annotator acts as an oracle that corrects the
27 model predictions and supplies sensitive information when necessary. This approach
28 helps to reduce the corpus construction time since annotators can focus on simpler
29 pieces of annotation.

30 Even with computer assistance, the amount of human effort, the required expertise,
31 and the time spent to build corpora, make this task an expensive one. Automatic
32 corpus extension through bootstrapping techniques is an effective way to increment
33 the corpus creation efficiency [11, 12]. Corpus bootstrapping denotes the process of
34 (1) manually annotating a small subset of a corpus, (2) training a model to predict
35 the annotations, and (3) using that model to automatically annotate the remaining
36 collection. This process can be repeated incrementally as long as the required quality
37 is ensured by some validation estimator. Even with low-performance systems, this
38 strategy is valuable if error correction is faster than manual annotation.

39 As with any automatic processing technique, there is a trade-off between quality
40 and efficiency that needs to be balanced. On the one hand, manual expert annotation
41 of the corpus should achieve the highest quality at the expense of efficiency. On the
42 other hand, automatic annotation allows a quicker annotation, but is sensitive to
43 overfit the initial collection and not generalize well to unseen situations. The use of

44 fine-tuned machine learning models can help to address this issue, but there is no
45 out-of-the-box solution.

46 1.1. Ensemble Methods

47 Ensemble methods are designed to address the low bias / high variance problem
48 exhibited by most machine learning models, making them suitable for producing more
49 robust classifications [13]. An ensemble model is made of several low bias models
50 whose predictions are combined to produce a final prediction. The main assumption
51 is that the combination of the low-level predictions will produce an output with a
52 lower variance while keeping a low bias. Having a diverse set of low-level models is a
53 key feature to achieve this [14].

54 The multi-hypothesis nature of ensemble models ensures that, if fine-tuned, it will
55 perform better than any of the individual models in the general case. This also allows
56 them to estimate the degree of confidence or quality of the predictions they output.
57 Classic ensembling techniques include *voting* and *weighted-voting* [15], *boosting* [16],
58 and *bagging* [17].

59 In a classification problem, *voting* produces as output the label that achieved the
60 majority of votes, treating each low-level model prediction as a vote. *Weighted-voting*
61 works just like voting, but each low-level model is assigned a weight that indicates
62 the importance of its votes. The label with the highest cumulative score is returned
63 as output. *Boosting* runs an iterative process where models are trained sequentially,
64 each one trying to improve its performance in the training examples that the previous
65 models performed the worst. During this process, each sub-model is also assigned a
66 score that weights the importance of its prediction. *Bagging* trains each sub-model
67 in a different selection (with replacement) of the original training examples. This
68 way, a model with a high variance should produce trained models with high diversity.

Reference annotation: "Cancer is a group of diseases"
Output annotation: "Cancer is a group of diseases"

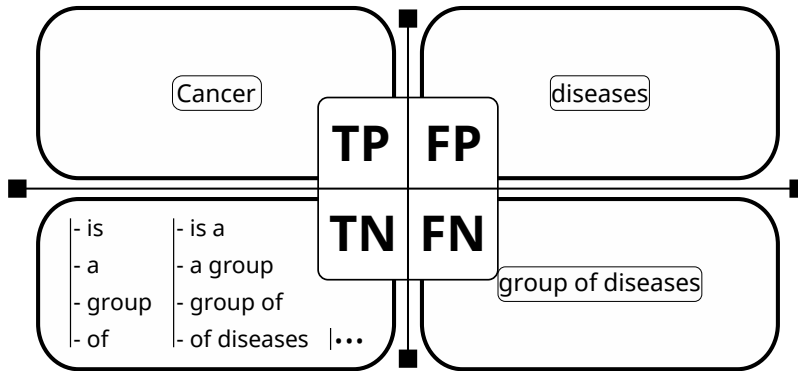


Figure 1: Classification of the annotations generated by a NER model according to a reference annotated version. **TP**, **TN**, **FP**, and **FN** stand for True Positive, True Negative, False Positive, and False Negative, respectively.

69 Alternatively, *feature bagging* works similarly by selecting a subset of the features
70 instead of the training examples, causing correlated features to be analyzed separately
71 in some sub-models.

72 While this works well for standard classification problems [18, 19], more complex
73 corpora require additional factors to be taken into account. For example, information
74 extraction tasks such as NER require a wider definition of what is being “voted” by
75 each sub-model and what should be reported as output. The NER task consists of
76 locating entities within the text by extracting spans of text and then assigning a
77 label to them. Overlapping entities with different labels produced by two sub-model
78 might imply that either the two entities are returned, or only one of them, or even
79 neither of them. Figure 1 shows how the annotations generated by a NER model are
80 classified into four categories according to a reference annotated version.

81 Some models might have better accuracy at predicting a certain type of entity

82 than others, and therefore, the weighted-voting should not be done uniformly across
83 all the different types. Advances in this area suggest that models should be weighted
84 independently for each entity class in order to obtain a better performance [20].

85 *1.2. Scenario for Ensembling*

86 Many information extraction competitions (e.g., eHealth-KD¹, IberLegal², and
87 CAPITEL³, at IberLEF 2020⁴) are hosted every year to encourage the research
88 community to produce new state-of-the-art results in the area. These competitions
89 allow the coordination of different systems to solve the same problem by defining
90 an evaluation corpus and a set of input and output specifications shared by all
91 systems. This turns each competition into a source of different annotated versions of
92 their corpora. By ensembling these results, the original corpus of the competition
93 can be extended and used as additional training resources in future versions of the
94 competition, which may have a positive impact on the performance of the participating
95 systems.

96 The eHealth-KD 2019 (or eHealth-KD v2 [21]) corpus, is a noteworthy example
97 of information extraction related corpora available in the literature [22]. It is made
98 of three types of annotations: entities, relations, and attributes, all of them manually
99 annotated in a collection of Spanish sentences from the medical domain. The corpus
100 was used as the main evaluation scenario for the eHealth Knowledge Discovery
101 Challenge at IberLEF 2019 (eHealth-KD 2019 for short) [23]. In the challenge,
102 participants were asked to automatically annotate a set of plain text sentences with

¹<https://knowledge-learning.github.io/ehealthkd-2020/>

²<https://temu.bsc.es/iberlegal/>

³<https://sites.google.com/view/capitel2020>

⁴<https://sites.google.com/view/iberlef2020/home>

103 the corresponding annotation scheme (attributes were excluded from the scope of
104 the challenge). Different systems were submitted and ranked according to their
105 performance in a gold standard collection. The participants were provided with a set
106 of 8 700 additional untagged sentences that were merged with the actual test collection.
107 As a result, a collection with several different versions of the same automatically
108 annotated sentences was obtained. A critical ensembling of these versions can be
109 used to automatically extend the previous corpus with a higher performance than
110 simply using the winning model.

111 This paper’s goal is to design and validate an ensembling strategy to automatically
112 extend Knowledge Discovery corpora, particularly oriented toward NER and relation
113 extraction tasks, from the systems’ outputs. This type of corpora extension resembles
114 an improved way of doing corpus extension through bootstrapping: instead of using a
115 single model to annotate the unlabeled sentences of a corpus, the predictions (outputs)
116 of several models are ensembled to produce a single, more robust, annotated version.
117 In this line, the eHealth-KD 2019 challenge was taken as a case of study, resulting
118 in the automatic annotation of a collection of 8 000 sentences [24]. The results are
119 compared against both classical voting systems as well as automatic learning models
120 according to their performance in the original evaluation scenarios. An upper bound
121 of the performance that can be achieved by only merging the submissions of the
122 challenge is presented. Additionally, the ensembling pipeline presented in this paper
123 was used as a system in the 2020 edition of the eHealth-KD challenge [25] by one
124 of the participating teams [26]. The team designed several deep-learning models to
125 solve the challenge and ensembled them to produce a final submission.

126 The specific contributions of this research are as follows:

- 127 • The introduction of ensemble methods for corpora extension as a better al-

128 ternative to classic corpus bootstrapping. The main advantages are reducing
129 the bias towards a particular architecture or set of rules and providing richer
130 information on the quality of the final corpus.

- 131 • The definition of an ensembling pipeline that works directly on the system
132 outputs rather than on the systems themselves. The main advantage is allowing
133 the integration of any group of systems, regardless of their architecture, and
134 non-programmatic annotation sources, i.e., human annotations.
- 135 • The automatic extension of the *eHealth-KD v2* corpus [21], a manually annotated
136 corpus of Spanish language sentences in the health domain. The newly annotated
137 sentences are publicly available as the *eHealth-KD 2019 ensembled corpus* [24]
- 138 • Comparison of the performance of several ensembling techniques in a particular
139 corpus and the study of the additional quality metrics obtained from extend-
140 ing the corpus using ensembling methods rather than classical bootstrapping
141 methods.

142 The remainder of the paper is organized as follows. Section 2 provides relevant
143 details on the eHealth-KD challenge and its corpus. Section 3 presents the ensembling
144 pipeline proposed in this paper, along with the different ensembling techniques it
145 includes. Section 4 shows the results of evaluating the ensembler, along with the
146 most promising strategies found. Additionally, it provides some statistics about the
147 collection of sentences that resulted from ensembling the eHealth-KD challenge’s
148 submissions. Section 5 discusses the most relevant aspects of the whole research.
149 Finally, Section 6 presents the conclusions, along with future lines of work.

150 2. Related Work

151 This section provides further information on the eHealth-KD challenge and its
152 corpus. The methods presented in this paper were directly applied and evaluated in
153 the resources provided by the challenge.

154 2.1. eHealth-KD Challenge and Corpus

155 The eHealth-KD challenge hosted at IberLEF proposes a set of resources and
156 evaluation scenarios to encourage the development of systems for the automatic
157 extraction of knowledge from unstructured text. Both the relevant entities and
158 the relations between them that occur in plain text documents are asked to be
159 identified. The F_1 score is used as the main evaluation metric. Both partial matching
160 of annotations and label misclassification are considered.

161 Figure 2 shows an example annotation of three sentences. Each entity annotation
162 is identified by its span of text, i.e., the position that the entity holds in text. Each
163 relation annotation is identified by the pair of entities it relates. Four types of entities
164 are defined: Concept, Action, Predicate, and Reference. The annotation scheme has
165 also 13 types of relations. Entity and relation representations match the specification
166 described in Section 3.2.

167 The challenge had the participation of 10 different contestant systems. Three
168 evaluation scenarios were developed, each one oriented toward entity detection, relation
169 extraction, or both tasks. The *main scenario* (Scenario 1) asked participants to output
170 both the entities and relations present in a collection of plain text sentences. The
171 *entity detection task* (Scenario 2) asked participants to find and label only the relevant
172 entities present in a collection of sentences, and, therefore, relation annotations were
173 ignored. For the *relation extraction task* (Scenario 3), both the sentences and entities
174 were given to the participants for them to output the relations occurring between

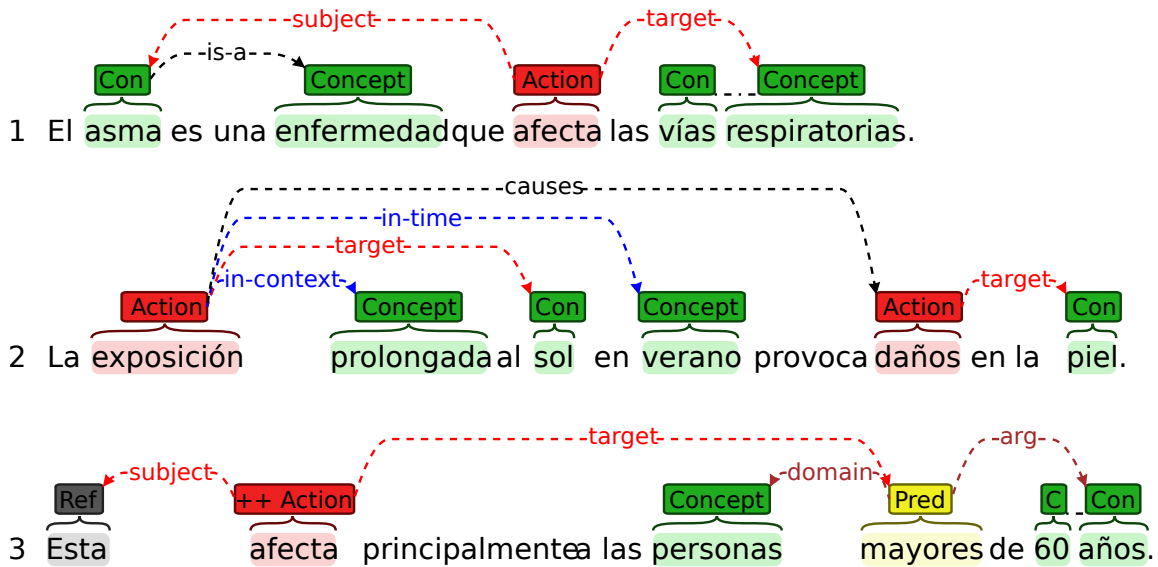


Figure 2: Example annotation of three Spanish sentences. The annotation shows the most relevant entities and relations defined.

175 them. A total of 32 successful submissions were made to Scenario 1, 48 to Scenario 2,
 176 and 48 to Scenario 3.

177 Each scenario is made of 100 sentences, for a total of 300 sentences with a
 178 human-annotated equivalent. A set of 8 700 additional sentences was added to the
 179 evaluation sentences of Scenario 1. Though these sentences were not used to measure
 180 the performance of the participant (since no reference annotation was available for
 181 them), this enforced participants to provide an automatically annotated version of
 182 the sentences while also discouraging them to manually annotate the test collection
 183 to cheat the evaluation system.

184 The ensembling pipeline presented in this paper can be applied to make use of the
 185 different versions⁵ of the 8 700 non-gold-annotated sentences to automatically extend

⁵<https://github.com/knowledge-learning/ehealthkd-2019/tree/master/data/>

186 the corpus. These versions consist of the submissions made by the contestants to the
187 additional sentences included in the test collection.

188 **3. Method**

189 This section presents our ensembling pipeline (also referred to as “the ensembler”),
190 its input, the set of algorithms and parameters it has available to select from, the
191 optimization algorithm used to find the best configuration, and an upper bound of
192 ensemble-based method’s performance. Section 3.1 provides a general overview of the
193 ensembling pipeline. Section 3.2 describes the input expected by the ensembler, along
194 with the requirements imposed on it. Section 3.3 presents the pool of algorithms that
195 the ensembler has available to select from, along with their parameters. Section 3.4
196 describes how to build a fake annotation system, called the “Oracle Ensembler”, to
197 obtain an upper bound of the performance obtained by ensembler systems based
198 only on other systems’ outputs. Finally, Section 3.5 defines an agreement metric
199 to enhance the information about the collections that result from the ensembling
200 pipeline.

201 *3.1. Overview*

202 The ensembling pipeline takes as input a collection of different annotated versions
203 of a document and produces a single annotated version of it (d^e). A key feature of the
204 ensemble is its ability to work directly with different versions of a document rather
205 than with the program or individual who generated it. The particularities of the
206 input expected by the ensembler are described in Section 3.2.

submissions/all

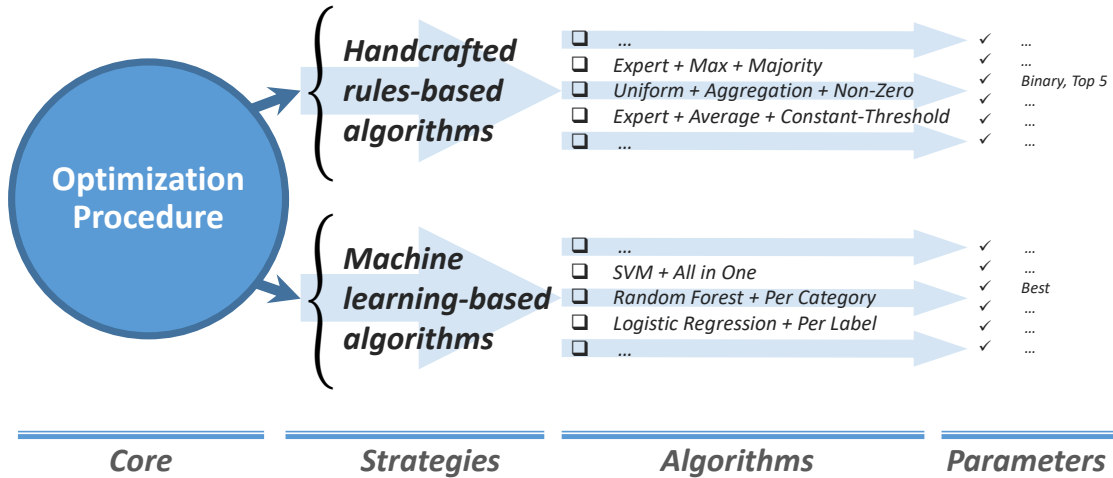


Figure 3: Overview of the ensembling pipeline. The ensembler’s core is an optimization procedure that searches for the best configuration according to a reference collection.

207 The ensembler consists of a pool of algorithms and parameters from which the
 208 “best configuration” is selected. The best configuration is defined as the combination
 209 of algorithms and parameters that generates the ensembled document with the
 210 highest score according to a predefined evaluation metric F and a reference annotated
 211 document (d^*). Figure 3 shows a high-level definition of the ensembler.

212 The F_1 metric is the standard fitness function to be used due to the information
 213 extraction nature of the NER and relation extraction tasks. However, any metric
 214 that scores the quality of a document d^e with respect to the reference d^* can be
 215 used. For example, additional considerations, such as partial match between entities,
 216 can be taken into account depending on the use case. A macro- F_1 metric can be
 217 applied to give each annotation type the same weight (even when some of them are

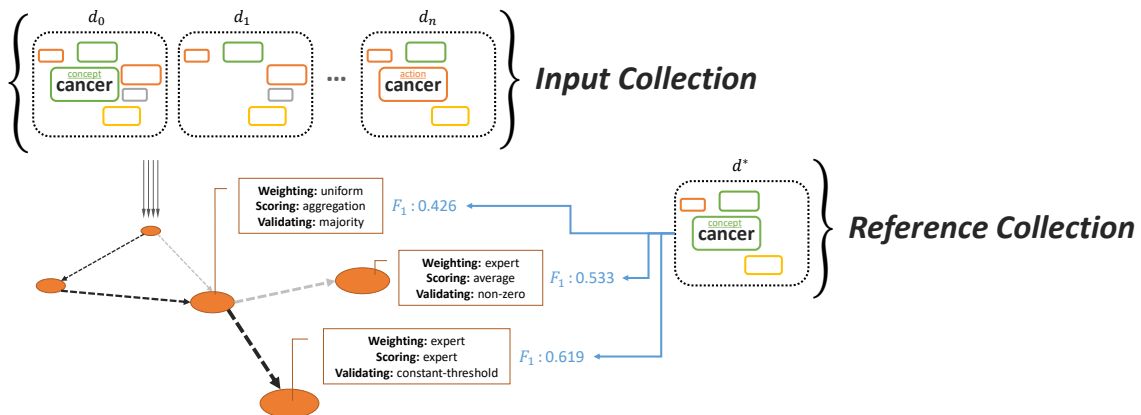


Figure 4: Optimization search performed on top of the pipeline configuration space. Each d_i , with $i \in [0, n]$, stands for an annotated version of the documents; d^* stands for the reference version of the document.

218 less represented in the whole collection). To apply these metrics, annotations are
 219 categorized into four classes, as shown in Figure 1. Precision is computed according
 220 to the number of True Positives and False Positives, which are the number of correct
 221 and spurious annotations, respectively. Recall is computed according to the number
 222 of True Positives and False Negatives, which are the number of correct and missing
 223 annotations, respectively.

224 To find the best configuration, a probabilistic evolutionary search is performed on
 225 top of the space of possible configurations (algorithms and parameters). The search
 226 starts with a random sampling strategy, but as it evaluates more pipelines, it modifies
 227 a probabilistic sampling model so that pipelines similar to the best ones found are
 228 more commonly sampled. Figure 4 illustrates the optimization procedure. The fitness
 229 function is set to directly maximize the F score of the resulting ensemble according
 230 to the reference collection.

231 The ensembler’s pool of algorithms is divided into two categories: (1) the hand-
 232 crafted rules-based algorithms and (2) the machine learning-based algorithms. The

233 first group contains the algorithms that apply a sequence of handcrafted rules to merge
234 the different versions of a document into a single ensembled version. The other group
235 contains the algorithms that train one or more machine learning models to predict,
236 based on the annotations made in the input collection of versions, which annotations
237 to include in the ensembled version. Section 3.3 provides the full specification of the
238 algorithms that are considered by the ensembler.

239 3.2. Input

240 Given a set of different annotated versions of a document $C = \{d_0, \dots, d_n\}$, $n > 1$,
241 the goal of our ensembling system is to produce a single annotated version d^e of the
242 document. Each d_i , with $i \in [0, n]$, is an annotated version of the document, also
243 referred to in this paper as a system, in reference to a system that generated the
244 document as output. An annotated document, in the scope of this paper, consists of a
245 set of annotated sentences $\{s_0, \dots, s_p\}$. The ensembler explores a space of parameters
246 in search of the configuration that produces the best ensembled document according
247 to a reference collection d^* . That is, the reference annotated version of some sentences
248 is provided, and their automatically ensembled versions are evaluated with respect to
249 them using a fitness function F .

250 Only entity and relation annotations are considered by our ensembler. An entity is
251 represented by a pair $\langle spans, type \rangle$, where $spans$ is a sequence $\{\langle b_0, e_0 \rangle, \dots, \langle b_n, e_n \rangle\}$
252 of spans of text in the sentence, with each $\langle b_i, e_i \rangle$ representing the beginning and end
253 of the i -th span, and $type$, a label that indicates the type of entity it was assigned.
254 Each entity belongs to a single sentence, i.e., all their spans of text are contained
255 in the same sentence. There might be overlapping between entities spans, but it is
256 assumed that there is no full overlap between entities of the same type, i.e., there is
257 no pair of entities with the same type that share the exact same span of text. Figure 5

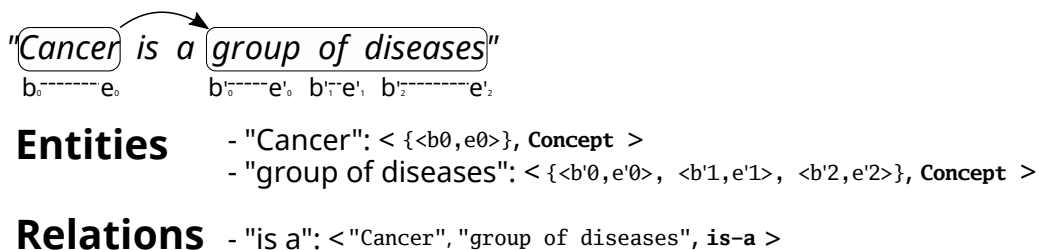


Figure 5: Summary of the representation of entities and relations expected by the ensembler.

258 shows a sample sentence with annotated entities and their representation.

259 A relation is represented by a triplet $\langle E_s, E_h, type \rangle$, where E_s and E_h are the
 260 source and head entities of the relation respectively, and $type$, a label that indicates
 261 the type of relation annotated. Relations occur only between entities that belong
 262 to the same sentence, i.e., there are no edges between entities in different sentences.

263 Figure 5 shows a sample sentence with an annotated relation and its representation.

264 3.3. Algorithms and Parameters

265 Two different types of algorithms are explored by the ensembler's optimization
 266 procedure as presented in Section 3.1:

- 267 • Those that pipeline a set of handcrafted rules to combine the collection of
 268 documents C (as described in Section 3.3.1).
- 269 • Those that train machine learning models to predict which annotations from
 270 the collection of documents C to include (as described in Section 3.3.2).

271 3.3.1. Handcrafted Rules-based Strategy

272 The ensembling pipeline has available a set of handcrafted rules-based algorithms
 273 to ensemble the collection of documents C . These algorithms sequentially apply a
 274 combination of manually designed rules to merge the output of several information

275 extraction systems. First, the ensemble determines, for each annotation in the union
276 of all input documents, which systems include it and which ones do not; we say that
277 systems “vote” for the annotations included in their output. Secondly, in the case
278 of conflicting annotations, the prevalent one(s) must be selected. Afterwards, the
279 ensembler must decide which proposed annotations are good enough to be included.
280 All of these decisions have an important role in producing the final output.

281 According to this, handcrafted rules-based algorithms, grouped into the next
282 Strategy Development section, are organized as follows:

- 283 1. **Vote Aggregation:** Given the collection of annotated documents C , the union
284 set of all annotations is built while keeping track of which system voted for each
285 annotation. It is said that system i voted for an annotation r if r is included in
286 document d_i . More details below in Section 3.3.1.
- 287 2. **Vote Weighting:** The votes given by each system to each annotation are
288 weighted. More details below in Section 3.3.1.
- 289 3. **Annotation Scoring:** An aggregated score of the quality of each annotation
290 is computed according to its votes. More details below in Section 3.3.1.
- 291 4. **Validation and Selection:** A decision is made as to whether the annotation
292 is kept and, in the case of conflicting annotations, a decision is reached on which
293 one to keep. More details below in Section 3.3.1.

294 Figure 6 summarizes the previous architecture. Several ensembling strategies
295 can be generated by assigning algorithms to all previous steps. For example, a
296 classical voting ensembler assigns a uniform weight to each voting system. Only the
297 annotations that exceed a predefined count are reported and in the case of conflict,
298 only those that achieve the majority of votes. Another example is an expert ensembler,

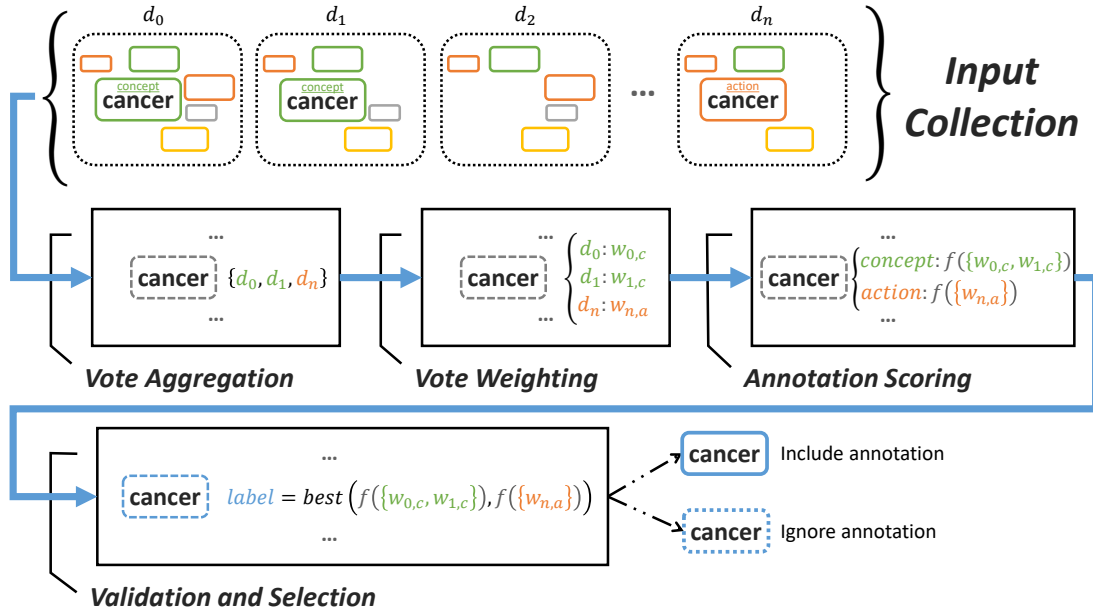


Figure 6: Summary of the handcrafted rules-based strategy. All annotations are aggregated and, afterward, a pipeline comprising three tasks decides which annotations to include in the ensembled collection. Each d_i , with $i \in [0, n]$, stands for an annotated version of the documents; each $w_{i,j}$ stands for the weight assigned to system i at predicting annotations of type j ; f stands for a function that scores the quality of annotation type j given all the active weights for type j ; $best$ stands for a function that selects the annotation type that was scored the highest (arbitrarily chosen in the case of a tie) and decides whether it is high enough to be included in the output collection.

299 which can be implemented by precomputing the importance per annotation type of
 300 each system and using it to weigh the votes. The votes emitted by the corresponding
 301 best system are weighted accordingly, and others valued as zero. Only the annotations
 302 with a non-zero vote are reported and in the case of conflict, only those that achieved
 303 the best score.

304 *Strategy Development*

305 The next sections provide details on each stage that takes part in the pipeline of
306 handcrafted rules-based algorithms.

307 ***Vote Aggregation.*** The ensembler has two ways to aggregate the annotations:

- 308 • *binary decisions*, in which the ensembler considers each annotation with a fixed
309 label as a possible candidate.
- 310 • *non-binary decisions*, in which the ensembler has to select only one label from
311 all the proposed labels for the same annotation.

312 The vote aggregation phase has a critical impact on the final ensemble collection.
313 Depending on whether binary decisions are made, the output will have full overlapping
314 annotations or not. Binary ensembling can build collections more similar to the
315 reference collection than non-binary ensembling since the output of the latter is a
316 subset of the former. On the other hand, non-binary ensembling should be easier to
317 tune since the final output is relaxed across all versions, i.e., at most one annotation
318 will be reported per overlapping annotation, instead of multiple ones.

319 ***Vote Weighting.*** Votes made by each system are weighted in two ways:

- 320 • The *uniform voting* approach is the most basic strategy, in which every vote is
321 considered to have the same weight. This strategy, combined with others later
322 in the pipeline, can produce some classic ensembling algorithms. For example,
323 by adding up the weighted votes per label of each annotation, assigning to the
324 annotation the label that most systems voted for, and keeping it only if the
325 score was greater than 50% of the total number of systems, a classic voting
326 ensembler can be built.

327 • An *expert-based weighting* approach exploits the fact that some systems have
328 better performance than others (according to the reference collection). The
329 quality of each system’s votes is estimated using the fitness function F . Since
330 systems might perform differently at each type of annotation, the weighting
331 function depends not only on the system identifier but also on the type of
332 annotation that is being voted. Each weight $w_{i,j}$ is computed by measuring
333 the performance of system i according to function F while considering only the
334 annotations of type j in both the reference and submitted collections.

335 ***Annotation Scoring.*** Once the collection of weighted votes per label has been built
336 for each annotation, all candidate labels are assigned a score. This score is meant to
337 capture the degree of confidence given by the voting systems, indicating whether the
338 label should be included. Several strategies are implemented:

- 339 • *aggregation scoring*, in which all weighted votes are added up.
- 340 • *average scoring*, in which weighted votes are averaged.
- 341 • *max scoring*, in which the maximum weighted vote is returned.
- 342 • *expert scoring*, in which the weight of the highest weighted system is returned
343 if voted by the same, or zero otherwise.
- 344 • *union scoring*, in which the maximum score (1) is returned if at least one system
345 voted.

346 The scoring can be done relative to the voting systems only or to the whole
347 collection of systems. In the former case, the fact that a system did not vote does
348 not influence the resulting score, while that is not for the latter case. This is valid,

349 for example, for *average scoring*, in which added weighted votes can be divided either
350 by the number of voting systems or the total number of systems. Another possible
351 modification to the previous strategies is to consider only the top k best-weighted
352 votes, for any value k between 1 and the total number of systems.

353 ***Validation and Selection.*** Depending on whether the pipeline decided to use
354 binary or non-binary decisions for classification, each annotation will have a single
355 candidate label or a list of them, respectively. In the first case, annotations with
356 different labels are being treated as different annotations and therefore validated
357 independently. In the second case, overlapping annotations have been grouped, and
358 each candidate label has been scored. The selected label will be the one that was
359 scored the highest (arbitrarily chosen in the case of a tie), resulting in a single
360 annotation per group of overlapping annotations. In either case, the final candidate
361 label has to be validated to decide whether it is good enough to be included in the
362 ensembled collection.

363 Different strategies are used:

- 364 • *non-zero validation*, in which an annotation will be reported if the label selected
365 for it has a positive score.
- 366 • *majority validation*, in which an annotation will be reported if the label selected
367 for it has a score greater than half the total number of systems.
- 368 • *per-label-threshold validation*, in which the score for a label has to exceed a
369 predefined (per label) threshold for an annotation with that label to be reported.
- 370 • *constant-threshold validation*, which works similarly to the *per-label-threshold*
371 *validation* but uses the same threshold for every type of label.

372 All the previous strategies are particular cases of the *per-label-threshold validation*,
373 in which the thresholds are set to zero for *non-zero validation*, to half the number of
374 systems for *majority validation*, or to parametric constant value for *constant-threshold*
375 *validation*. Defining these particular cases is useful to improve the performance while
376 searching the space of possible configurations, as described in Section 3.1. Some
377 of these validations are oriented toward certain strategies previously applied. For
378 example, the *majority validation* intuitively should have a better performance when
379 combined with *uniform weighting* and *aggregation scoring*, since that represents a
380 classic voting ensembler.

381 *Optimization Parameters*

382 The task of deciding which strategies and parameters combine the best is given
383 to the optimization search algorithm as presented in Section 3.1.

384 Table 1 summarizes the space of handcrafted-rules based algorithms to be explored.
385 Strategies are chosen from categorical values. Their parameters are requested to the
386 space sampler, which is then updated according to the performance achieved by the
387 pipeline.

388 In some case, the input systems are highly correlated or might belong to the same
389 case study. For example, in the case of the eHealth-KD challenge, each participant
390 might submit multiple systems, which in some cases share the exact same architecture
391 but with small fine-tuned parameter changes. Having multiple systems biased towards
392 the same type of annotations might (or might not) have a bad influence on the final
393 ensemble. To deal with this, the *best* parameter was added to the optimization space.
394 This parameter controls whether to consider only the best performing system of each
395 category (e.g., the best performing submission of each participant) or all of them. The
396 inclusion of the *best* parameter in the building phase allows the optimization algorithm

Phase	Strategy	Parameters
	voting	binary ^b
Building		best ^b
	submissions	top ^k
		selection ^l
Weighting	uniform	
	expert	
Scoring	aggregation	top ^k
	average	top ^k , strict ^b
	max	
	expert	discrete ^b
	union	
Validating	non-zero	
	majority	
	constant-threshold	threshold ^t
	per-label-threshold	thresholds ^{t*}

Table 1: Summary of the handcrafted rules-based algorithms and their parameters to be explored. Each tag b , k , l , t and t^* stands for Boolean, integer, list, continuous, and vector type, respectively.

397 to turn it off if proved convenient. Additionally, the *top* and *selection* parameters
398 were added to ensemble the top best-performing systems only or a selection of them.

399 3.3.2. Machine Learning-based Strategy

400 Aside from the handcrafted rules, the ensembling pipeline has also available a
401 set of machine learning-based algorithms to ensemble the collection of documents C .
402 Machine learning models are trained to predict the probability of a labeled annotation
403 to be correct, i.e., the probability of it belonging to a gold annotated collection,
404 according to the votes given to it by each system. Two key aspects are configured
405 here: first, how each annotation is transformed into a feature vector and which
406 model handles it, and second, which concrete model architecture is used. Figure 7
407 summarizes the general structure of this strategy.

408 *Strategy Development*

409 **Annotation Representation** The main features for representing any labeled
410 annotation comes from which systems voted for it. That means that part of the
411 feature vector will hold the information of whether a system voted or not for the
412 labeled annotation to be included. This is captured with a vector $v = \langle v_0 \cdots v_n \rangle$, n
413 being the total number of systems, in which $v_i = 1$ if system i voted to include the
414 annotation with that label, and $v_i = 0$ otherwise. Instead of using binary weighting,
415 the importance of the corresponding voting system can also be used, this is, $v_i = w_{i,j}$
416 if system i voted to include the annotation with label j , with $w_{i,j}$ being the system
417 weight. Additionally, a one-hot encoded vector l is built to encode the selected label,
418 this is, l is a vector whose size equals the number of possible labels, and in which all
419 component equals 0 but the one that corresponds to the selected label. Using this
420 information, the feature vectors and models can be handled in one of the following
421 ways:

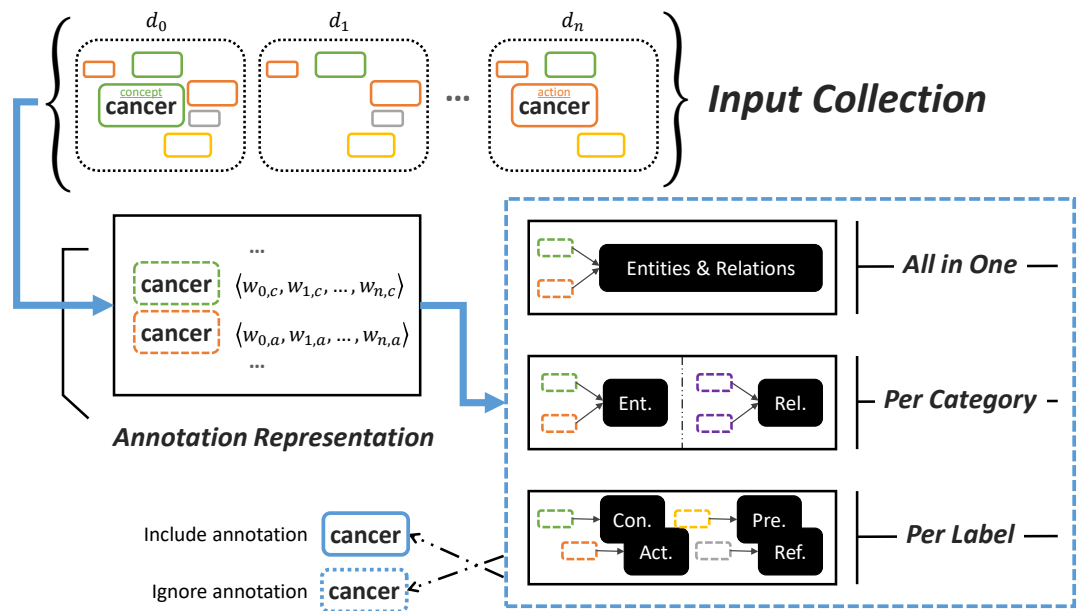


Figure 7: Summary of the machine learning-based architecture. Each annotation is transformed into a feature vector and, afterward, a machine learning model computes the probability of it belonging to a gold annotated collection. Each d_i , with $i \in [0, n]$, stands for an annotated version of the documents; each $w_{i,j}$ stands for the weight assigned to system i at predicting annotations of type j .

- 422 • *all in one model*, in which a single model is trained to handle every labeled
423 annotation. Each feature vector is built by concatenating vector v and vector l .
- 424 • *per category model*, in which two models are trained, one to handle entities and
425 one to handle relations. Each feature vector is built by concatenation vector v
426 with vector l (using the entity or relation labels accordingly).
- 427 • *per label model*, in which a model is trained for each type of label. Each feature
428 vector is fully composed of vector v .

429 **Model Architecture** In the context of this paper, three classic *scikit-learn*⁶
430 classifiers were considered: Random Forest, SVM, and Logistic Regression. The
431 default parameters of the respective models in the 0.22.2.post1 version of the
432 library are used. The three architectures were evaluated as reported in Section 4. In
433 general, any machine learning classifier can be registered to solve the prediction task.
434 The task consists in a binary classification problem in which the feature vector of
435 each labeled annotation is used to predict whether to output the annotation or not.

436 *Optimization Parameters*

437 Table 2 summarizes the space of machine learning-based algorithms to be explored
438 by the optimization procedure described in Section 3.1. Aside from the weighting
439 and training modes and the model architecture, an additional parameter “best” was
440 included. This parameter controls whether to use all systems or only the best one
441 per category. Binary voting is used in this situation for vote representation to fit the
442 learning systems requirements.

⁶<https://scikit-learn.org/>

Phase	Strategy	Parameters
Building	submissions	best ^{<i>b</i>}
Weighting	uniform	
	expert	
Model	SVM	
	Random Forest	
	Logistic Regression	
Training	all in one	
	per category	
	per label	

Table 2: Summary of the space of machine learning-based algorithms to be explored. Tag *b* stands for boolean.

443 3.4. Upper Bound

444 Though ensemble methods are a great way to improve the performance of individual
445 systems, the diversity of those systems plays a key role in the quality of the final
446 results. There is a limit to what an ensemble method can achieve based on its
447 subsystem. To obtain an upper bound of the performance that the best ensembling
448 method can achieve, an oracle ensembler can be built.

449 The oracle ensembler consists in an algorithm that always makes the correct choice
450 on which voter to trust. That means that for a given annotation, even if only one
451 system voted for the inclusion of the annotation, or only one system voted for the
452 exclusion of the annotation, the oracle ensembler will trust it if doing so is the right
453 decision according to the reference collection. As a consequence, every annotation
454 present in the reference collection that was identified by at least one system will be
455 included. On the other hand, annotations that were not detected by any subsystem
456 will not be included. Only if all subsystems vote for the inclusion of an annotation
457 not present in the reference collection, a spurious annotation will be reported.

458 It only makes sense to build such an ensemble for comparison purposes, since
459 it requires a gold annotated sentence to make a prediction on that same sentence,
460 which means that it will solve a problem already solved. Section 4 shows the obtained
461 upper bound of the performance of any ensembling system based only on the previous
462 submissions of the eHealth-KD challenge can obtain.

463 3.5. Quality Measures

464 The use of ensembling methods to produce a single merged version from different
465 versions of a document allows computing an agreement metric to further enhance the
466 information available on the quality of the final result. Though this condition is not
467 determinant in the quality of the resulting ensemble (since part of the performance

468 of ensemble methods comes from the diversity of their subsystems), the agreement
469 between annotations estimates the confidence in the final document. Sentences that
470 have a higher agreement are probably easier to annotate. That means that a corpus
471 made of those sentences is more likely to be correctly annotated since subsystems
472 may have achieved better performance on their own.

473 To measure the agreement in a sentence, the following equation is proposed:

$$agreement(anns) = \frac{\sum_{i \in anns} |votes(i)|}{N \cdot |anns|}$$

474 where *anns* is the set of annotations in the sentence, *votes* returns the systems than
475 voted for the annotation, and *N* is the total number of systems. This metric value
476 ranges from 0 to 1. The highest value is achieved when every annotation received the
477 maximum number of votes (*N*).

478 4. Results

479 The methods described in this paper were directly applied and evaluated in the
480 resources provided by the eHealth-KD challenge. The eHealth-KD challenge provides
481 a collection of 9 000 sentences automatically annotated by each participant system
482 of the challenge’s 2019 edition.⁷ The gold-annotated version (manually annotated
483 by experts) of 600 of those sentences is also available. Those manually annotated
484 sentences are divided into two packs of 300 sentences, used to evaluate the submissions

⁷<https://github.com/knowledge-learning/ehealthkd-2019/tree/master/data/submissions/all>

485 of the challenge’s 2019⁸ and 2020⁹ editions, respectively.

486 In the context of this paper, the previous collections of sentences are used in the
487 following way:

- 488 • **Reference Collection [300]** The sentences used to evaluate the submissions
489 of the challenge’s 2020 edition. Both the gold-annotated version and the
490 submissions of each participant (to the challenge’s 2019 edition) are available.
491 Conveniently, these sentences were selected from different agreement ranges
492 (as described in Section 3.5), which means that some of them (100) had high
493 agreement, others (100) low agreement, and the remaining (100) an in-between
494 agreement. [27]
- 495 • **Validation Collection [100]** The sentences used to evaluate the submissions
496 to the main scenario (Scenario 1) of the challenge’s 2019 edition. Both the
497 gold-annotated version and the submissions of each participant are available.
- 498 • **Blank Collection [8400]** The remaining sentences, for which the gold-
499 annotated version is not provided. Only the submissions of each participant to
500 the challenge’s 2019 edition are available.

501 Three different aspects of the ensembler described in this paper were evaluated.

⁸Available in the testing collection (<https://github.com/knowledge-learning/ehealthkd-2019/tree/master/data/testing>) of the eHealth-KD 2019 challenge (scenarios 1, 2, and 3).

⁹Available in the testing collection (<https://github.com/knowledge-learning/ehealthkd-2020/tree/master/data/testing>) of the eHealth-KD 2020 challenge (scenarios 1, 2, and 3).

- 502 • **Quality of the Ensembled Corpus** First, the ensembler was used to produce
503 the *eHealth-KD 2019 ensembled corpus* [24], a collection of 8 000 automatically
504 annotated sentences, from the Blank Collection’s sentences. Of these, 3 000
505 were published in the context of the eHealth-KD 2020 challenge as an additional
506 source of development data¹⁰. The Reference Collection was used as a reference
507 collection for tuning the ensembling pipeline. Section 4.1 presents some statistics
508 about the final corpus, including an analysis of quality metrics.
- 509 • **Quality of the Ensembling Algorithms** Secondly, the performance of
510 the most interesting ensembling strategies is evaluated in the Reference and
511 Validation Collections. Section 4.2 presents these results in a greater detail.
- 512 • **Quality of the Ensembler Optimization** Finally, the performance of the
513 ensembler’s optimization procedure is evaluated. With this goal in mind, the
514 ensembler was used as a system to participate in the eHealth-KD 2020 chal-
515 lenge [25]. Team *UH-MatCom* ensembled several deep-learning architectures to
516 produce an additional submission to the challenge [26]. The ensemble run ob-
517 tained the best results out of the three runs of the team, except for an incorrectly
518 submitted scenario. This shows that the ensembling pipeline produces outputs
519 that perform better than the systems given as input. Section 4.3 elaborates on
520 these results.

521 4.1. Quality of the Ensembled Corpus

522 In this section, the experiments applied to evaluate the ensembler’s quality in
523 extending corpora are presented. The ensembler was used to produce the *eHealth-KD*
524 *2019 ensembled corpus* [24], a collection of 8 000 automatically annotated sentences [24],

¹⁰<https://github.com/knowledge-learning/ehealthkd-2020/tree/master/data/ensemble>

525 from the Blank Collection of sentences. The selection criteria was to keep the top
526 8 000 sentences with the highest agreement (as described in Section 3.5). Figures 8,
527 9, 10, 11, and 12 summarize the steps applied to produce the final collection of
528 automatically annotated sentences. The Reference Collection was used as a reference
529 collection for tuning the ensembling pipeline. The impact of using the ensembled
530 corpus as an extra source of training data for solving the challenge was tested, as
531 shown in Section 4.1.3. Additionally, the performance of already tuned ensembler was
532 evaluated in the three original evaluation scenarios of the eHealth-KD 2019 challenge,
533 as shown in Section 4.1.4. Finally, the quality of the resulting corpus was measured
534 with the agreement metric presented in Section 3.5, as shown in Section 4.1.5.

535 4.1.1. Ensembler Configuration

536 The ensembling pipeline was tuned in the 300 sentences of the Reference Collection.
537 This means the configuration was selected from the best configurations found in the
538 Reference Collection. Figure 8 illustrates the procedure.

539 Instead of the F_1 metric proposed in the challenge, a macro- F_1 variant of this
540 metric was applied in the ensembling pipeline. The macro- F_1 metric averages the F_1
541 obtained across each annotation type. This was done to ensure that all annotations
542 types were represented in the corpus. The F_1 metric used in the challenge benefits
543 systems that decide to ignore certain types of annotations to achieve a better score.
544 For example, the ensembler shown in Table 8 that resulted from optimizing the F_1
545 metrics decides not to report *causes*, *entails* and *has-property* relations. This is not
546 a problem of the ensembler but of the evaluation metric selected. The ensembler
547 was able to find a configuration that maximizes that metric. The correct metric to
548 be used is subject to the problem that is being solved. For example, in the context
549 of a system participating in the eHealth-KD challenge, maximizing the F_1 score is

ENSEMBLER CONFIGURATION

Reference Collection [300 sentences]

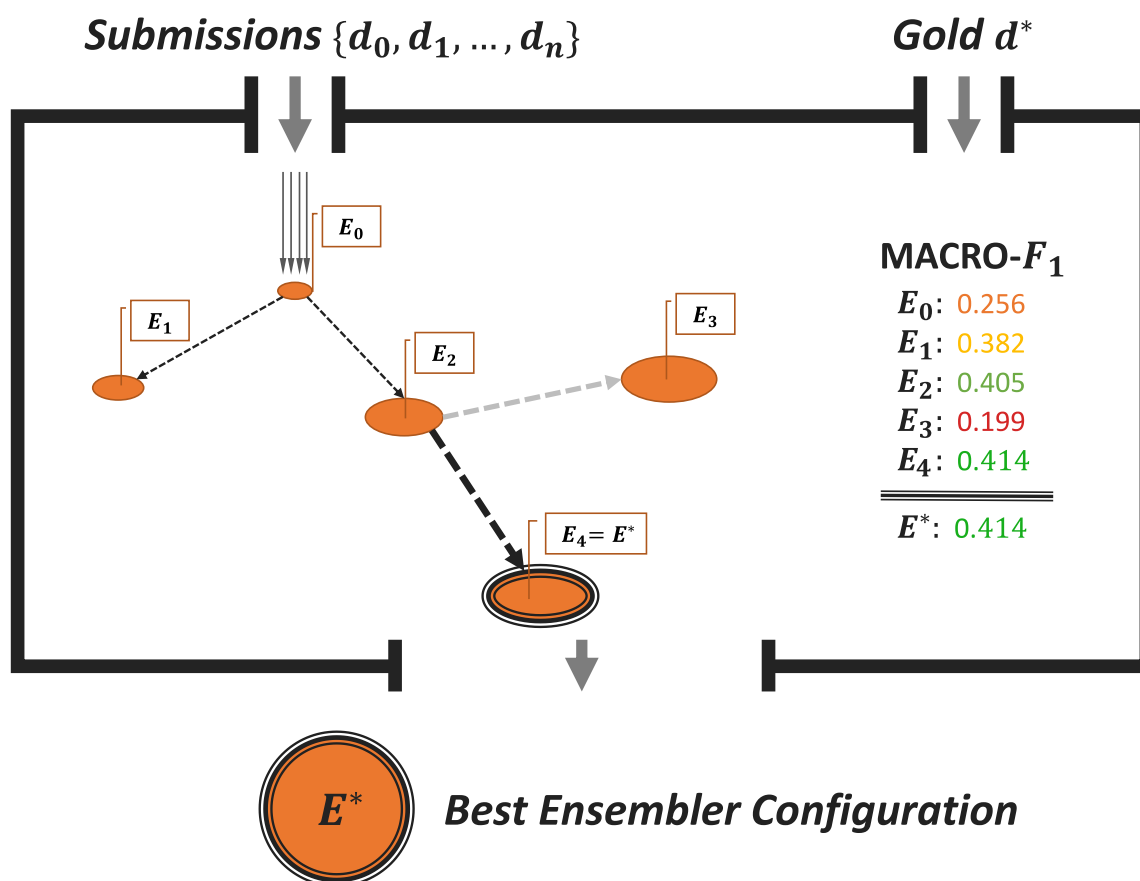
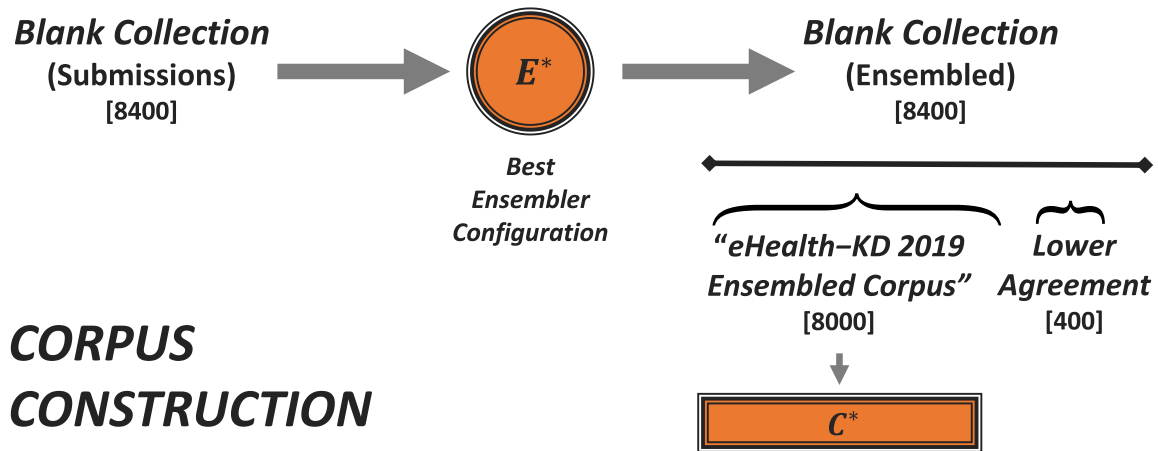


Figure 8: Ensembler configuration. The optimization search is performed using the Reference Collection. The submissions $\{d_0, d_1, \dots, d_n\}$ are ensembled using several configurations ($E_0 \dots E_4$) looking for the ensembler configuration E^* for which the macro- F_1 in the gold document is maximized. The found E^* configuration was used to build the *eHealth-KD 2019 ensembled corpus*. Section 4.1.1 provides further information on E^* .



CORPUS CONSTRUCTION

Figure 9: Corpus construction. The best ensembler configuration (E^*) found in the Reference Collection is used to build the *eHealth-KD 2019 ensembled corpus*. All sentences in the Black Collection are ensembled. From those 8 400 sentences, the 8 000 ones with higher agreement are included in the corpus (C^*). Section 4.1.2 provides statistics of the corpus (C^*).

550 the correct decision. However, for automatically annotated corpora, it is relevant to
 551 ensure that all annotation types are represented.

552 Table 3 summarizes the configuration that was used to ensemble the 8 000 sentences
 553 of the published *eHealth-KD 2019 ensembled corpus*. The macro- F_1 metric was used
 554 in the optimization search to measure the performance in the Reference Collection.
 555 The selected configuration achieved a macro- F_1 of 0.414 and an F_1 of 0.616 in the
 556 Reference Collection. Similarly, the selected configuration achieved a macro- F_1 of
 557 0.455 and an F_1 of 0.647 in the Validation Collection. Figure 9 illustrates the corpus
 558 construction procedure.

559 4.1.2. Corpus Statistics

560 Table 4 shows the number of entities and relations annotated in the *eHealth-KD*
 561 *2019 ensembled corpus*. A total of 86 112 elements were annotated, 50 249 of them
 562 as entities and the remaining 35 863 as relations. It is noteworthy that although

Phase	Strategy	Parameters
Building	voting	binary True
	submissions	best False
		top ∞ selection 15
Weighting	expert	
Scoring	expert	discrete True
Validating	constant-threshold	threshold 0.211

Table 3: Best ensembler configuration found according to the macro- F_1 metric. The ensembler configuration was used to produce the *eHealth-KD 2019 ensembled corpus*, a collection of 8000 automatically annotated sentences from the Black Collection. The Reference Collection was used as a reference collection for tuning the ensembling pipeline. Strategies and parameters are taken from Table 1. The collection of 15 selected eHealth-KD 2019 submissions is made of: `lsi2_uned/574805`, `lsi2_uned/575373`, `lsi_uned/575160`, `abravo/577194`, `iakesg/576597`, `iakesg/576606`, `vsp/576664`, `vsp/576666`, `vsp/576670`, `vsp/576671`, `jlcuad/577097`, `jlcuad/577119`, `talp/576647`, `hulat-taskAB/576454`, `hulat-taskAB/576998`.

563 attributes are part of the original annotation scheme of the *eHealth-KD v2* corpus,
564 they were not considered in the context of the challenge and, therefore, are not
565 available for the automatic construction of this ensembled corpus. Just as the original
566 corpus, the most represented entity classes are *Concept* and *Action* [22]. The most
567 represented relations are *target*, *subject*, and *in-context*.

Metric	eHealth-KD v2	Ensembled
Sentences	1 045	8 000
<i>Entities</i>	6 612	50 249
Concept	4 092	31 126
Action	1 742	13 522
Predicate	563	4 204
Reference	215	1 397
<i>Relations</i>	6 049	35 863
target	1 729	13 289
subject	894	5 144
in-context	677	4 975
is-a	566	3 284
in-place	400	1 580
causes	367	1 323
domain	364	1 217
argument	343	1 165
entails	167	423
in-time	165	565
has-property	159	2 242
same-as	124	467
part-of	94	189

Table 4: Summary statistics for the *eHealth-KD 2019 ensembled corpus*. The statistics of the *eHealth-KD v2* corpus are taken from (Piad et al. [22]).

CORPUS IMPACT ON A SYSTEM

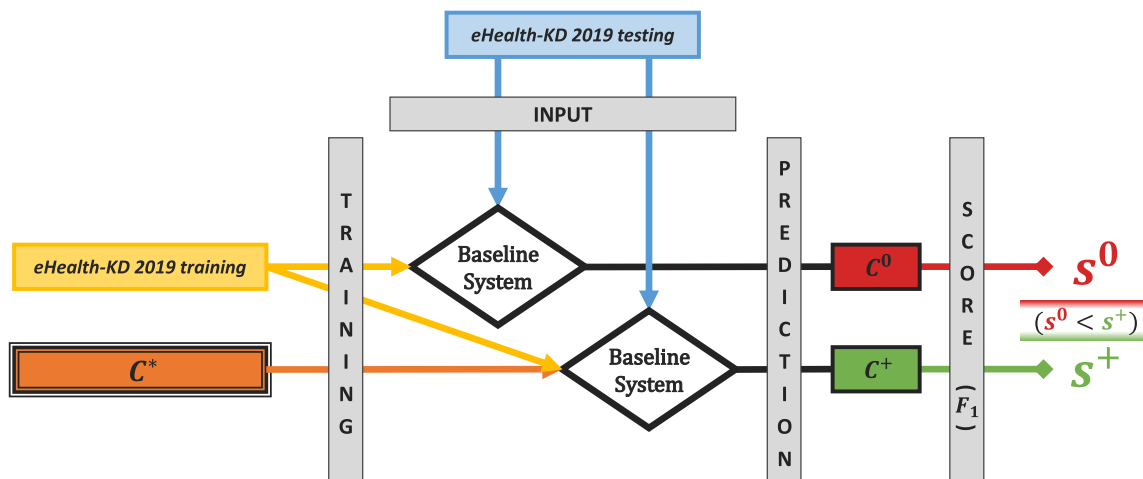


Figure 10: Corpus impact on a system. Two instances of the same system are trained: one using only the eHealth-KD 2019’s training collection, and the other one using also the *eHealth-KD 2019 ensembled corpus* (C^*). Both systems are then used to predict the eHealth-KD 2019’s testing collections, for which the F_1 score is computed as proposed in the challenge. The system that uses both training collection performs better than the other one, as shown in Section 4.1.3.

568 4.1.3. Corpus Impact on a System

569 To evaluate the usefulness of the 8 000 automatically annotated sentences as
 570 additional corpora, a system that used such collection as additional training data is
 571 evaluated. Specifically, the system was first trained using only the original training
 572 collection of the eHealth-KD 2019 challenge and another instance of the same system
 573 was then trained in a collection containing both the original sentences and the
 574 *eHealth-KD 2019 ensembled corpus*’s sentences. The original evaluation scenarios of
 575 the eHealth-KD 2019 challenge were used to evaluate both instances of the system.
 576 The F_1 metric proposed in the challenge was used in the three scenarios. Figure 10
 577 illustrates this procedure.

578 The selected system is a new baseline algorithm that was implemented and trained

579 to measure the impact of using the ensembled collection for future annotations.¹¹
580 The new baseline consists of two machine learning models, each one responsible for
581 solving one of the tasks of the challenge.

582 The new baseline solves the entity extraction task as a sequence labeling problem.
583 Sentences are tokenized and each token is assigned a label in the BILOUV [28] entity
584 tagging scheme¹². With the exception of the *Other* tag, all BILOUV tags are suffixed
585 with the four entity types, for a total of 21 labels. Several linguistic features (e.g.,
586 text, lemma, POS-tag, dependency tree label, etc.) obtained from *spaCy*¹³ are used
587 to represent the tokens. A CRF [29] estimator is used to predict the sequence of
588 output labels from the input sequence of tokens.

589 The new baseline solves the relation extraction task by using a Logistic Regression
590 model to predict which type of relation (if any) exists between each pair of entities in
591 the same sentence. The concatenation of the features of both tokens is used as input
592 to the model.

593 Table 5 summarizes performance achieved by the new baseline, when trained in the
594 original training collection only and when trained in both collections (the *eHealth-KD*
595 *2019*'s training collection plus the *eHealth-KD 2019 ensembler corpus*), for the three
596 evaluation scenarios. In every scenario, there is an improvement in performance
597 while using the information in the ensemble collection. This indicates that there
598 are new patterns not found in the original collection that can now be learned from
599 the ensembled collection. More complex systems might make a better use of this
600 information. A subset of the ensembled collection (3 000 sentences) is provided to the
601 participants of the eHealth-KD challenge at IberLEF 2020. Participants are free to

¹¹<https://github.com/jpconsuegra/ensemble-baseline>

¹²<https://devopedia.org/named-entity-recognition> (Accessed 2020-05-01).

¹³<https://spacy.io/>

Baseline	1-main	2-taskA	3-taskB
training	0.533	0.778	0.356
training + ensemble	0.556	0.800	0.363
human	0.727	0.861	0.735

Table 5: Comparison of performance of a baseline system when trained only on the original training collection provided in the eHealth-KD 2019 challenge, and that trained with the addition of the *eHealth-KD 2019 ensembled corpus*. The system is evaluated in the three evaluation scenarios of the eHealth-KD 2019 challenge. A human baseline is also provided for comparison purposes.

602 use this collection or not. The impact of using this additional information will be
603 later measured according to how the performance of their system was influenced by
604 the inclusion of this collection.

605 For reference purposes, the human baseline presented by Piad et al. [22] is
606 included in Table 5. The performance of the human baseline approximates the best
607 performance that can be achieved by any automatic annotation system. Relatively to
608 this score, the performance of the baseline system when trained only on the original
609 training collection provided in the eHealth-KD 2019 challenge, and that trained with
610 the addition of the ensemble collection, change to 0.733 and 0.765 in Scenario 1,
611 respectively. In Scenario 2, the scores change to 0.904 and 0.929, respectively. Finally,
612 in Scenario 3, the scores change to 0.484 and 0.494, respectively.

613 4.1.4. Estimated Quality

614 To estimate the quality of the selected ensembler configuration (shown in Table 3)
615 in annotating sentences and, therefore, the quality of the *eHealth-KD ensembled*
616 *corpus*' sentences, that same ensembler configuration was used to annotate a manually
617 annotated collection of sentences. The three evaluation collections of the eHealth-KD

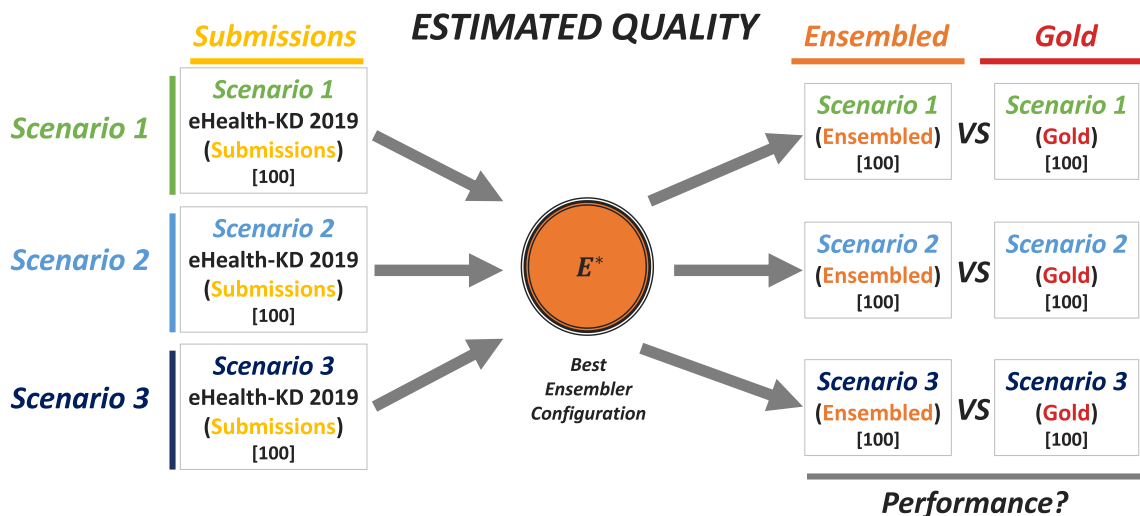


Figure 11: Estimated quality. The quality of the ensembler configuration E^* used to generate the *eHealth-KD 2019 ensembled corpus* is estimated by ensembling collections for which there are human-annotated versions. The three evaluation scenarios of the eHealth-KD 2019 challenge are used. Those collections were not used to tune the ensembler configuration E^* . Section 4.1.4 provides further information in this regard.

618 2019 challenge were used for this purpose. The sentences of those collections are
 619 not part of the Reference Collection used to tune the ensembler. Those sentences
 620 were manually annotated by the organizers of the challenge. Figure 11 illustrates the
 621 quality estimation procedure.

622 Table 6 compares in the three evaluation scenarios the performance achieved by
 623 the selected ensembler against the baseline and the best performing submission to
 624 the challenge (TALP-UPC). Additionally, the human baseline presented by Piad et
 625 al. [22] is included. The ensembler was able to outperform TALP-UPC’s performance
 626 in all scenarios but the third one. This scenario deals only with the extraction of
 627 relevant relations. It has been shown that this task is considerably harder than the
 628 entity recognition one [23]. The results achieved by the ensembler suggest that there

System	F_1			macro- F_1		
	1-main	2-taskA	3-taskB	1-main	2-taskA	3-taskB
baseline	0.431	0.546	0.123	0.182	0.461	0.095
TALP-UPC	0.639	0.820	0.627	0.451	0.792	0.520
ensemble (Opt.)	0.647	0.823	0.610	0.455	0.792	0.491
human baseline	0.727	0.861	0.735	0.592	0.830	0.666

Table 6: Summary of the performance, in the three evaluation scenarios, of the best ensembler configuration (according to the macro- F_1 metric) against the baseline, the best performing submission, and a human baseline.

629 is no fixed rule to follow for choosing which sub-system to trust. No model decided to
630 completely trust TALP-UPC’s annotations, which suggests that for certain situations
631 (present in the reference collection) trusting other systems might be a better option.
632 In the long run, this rule can not be extended to the whole collection of sentences.

633 The selected ensembler configuration can produce better annotations for Scenario
634 1¹⁴ than all of the individual systems. As shown in Table 6, this is true even for
635 sentences that were not present in the reference collection. This fact supports the
636 idea that corpora extension using our ensembling pipeline produces better results
637 than corpora extension through bootstrapping (since the latter would use only the
638 output of the best performing system, TALP-UPC in this case).

¹⁴Scenario 1 is considered the main evaluation scenario of the challenge and includes both tasks (the entity and relation extraction tasks).

639 4.1.5. Sentence Agreement

640 As explained in Section 3.5, the use of ensembling methods to produce a single
641 merged version from different versions of a document allows computing an agreement
642 metric to further enhance the information available on the quality of the final result.
643 In this section, the agreement on each *eHealth-KD ensembled corpus*' sentence is
644 computed. Then, it is shown that there is a relation between the agreement of the
645 sentences in a corpus and the quality of its annotations. Figure 12 illustrates this
646 procedure.

647 Figure 13 shows the number of sentences in the collection per agreement range.
648 Only one system per participant was considered to do this measure. This is done to
649 avoid biasing the results toward a specific participant if they submit a large number of
650 similar systems. Most sentences had an agreement around 0.2 and 0.3, which means
651 that in most sentences only 20 – 30% of the sub-systems agreed. This reveals how
652 difficult it is to get an agreement among many knowledge discovery systems.

653 Figure 14 shows the performance, of the ensembler configuration used to build the
654 *eHealth-KD ensembled corpus*, in different subsets of the Reference and Validation
655 Collections as the number of sentences increases, sorted by ascending (and descending)
656 agreement. As expected, as the number of sentences increases, the collections with
657 higher agreement tend to perform better than those with a lower agreement. This
658 shows that the agreement metric unlocked by the ensembler can be used to estimate
659 and control the quality of a collection of sentences for which there is no available
660 human-annotated version. The sentences published in the *eHealth-KD 2019 ensembled*
661 *corpus* [24] are sorted by agreement so that the top agreement sentences can be easily
662 filtered. The explicit agreement of each sentence is also provided.

Sentence Agreement

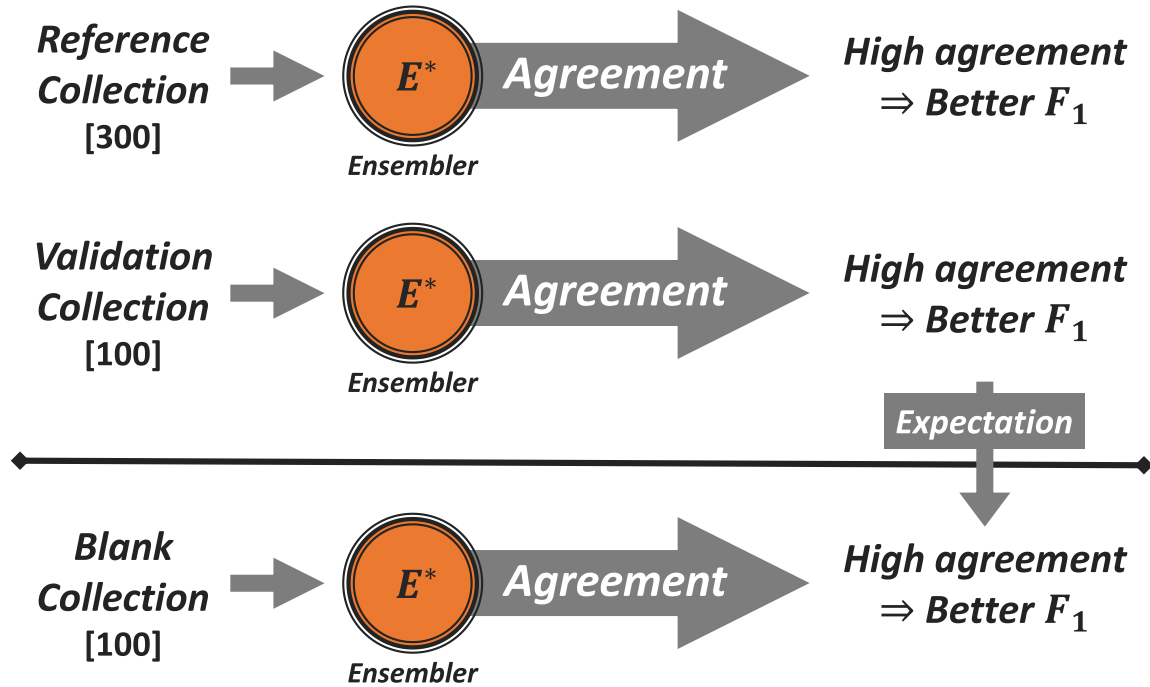


Figure 12: Sentence agreement. The agreement of all sentences in the Reference Collection and Validation Collection is computed. Such analysis shows that subsets of sentences with a higher agreement tend to have a higher F_1 score. This fact allows using the agreement of the sentences in the *eHealth-KD 2019 ensembled corpus* to control its quality. Section 4.1.5 provides further information in this regard.

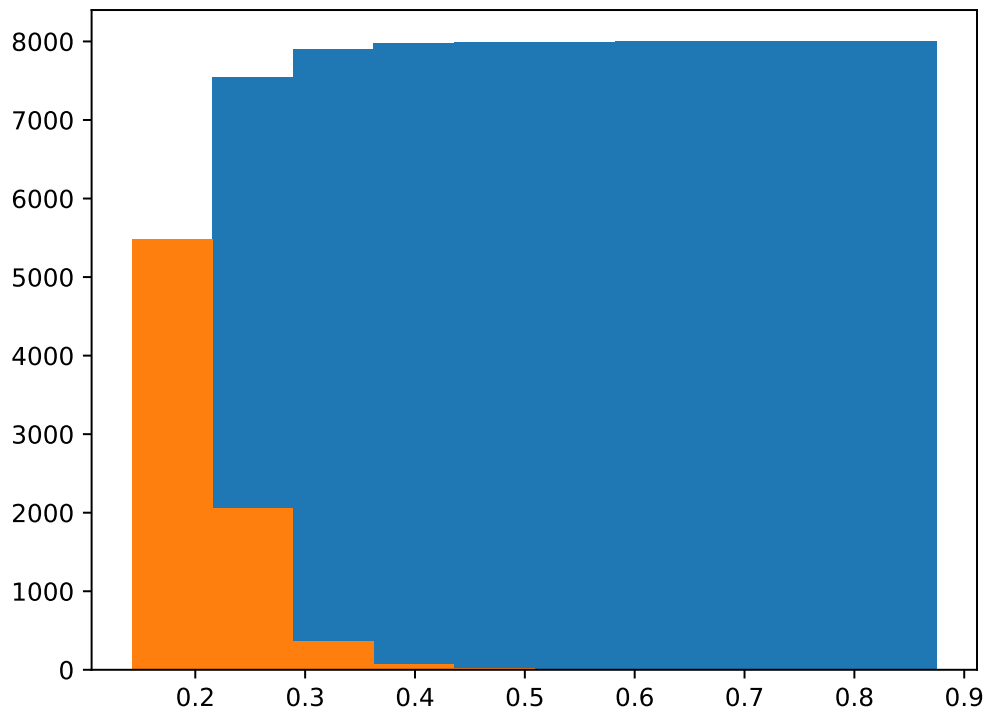


Figure 13: Number of sentences (in the *eHealth-KD 2019 ensembled corpus*) per agreement range between submissions. The histogram is shown in orange. The cumulative distribution is shown in blue.

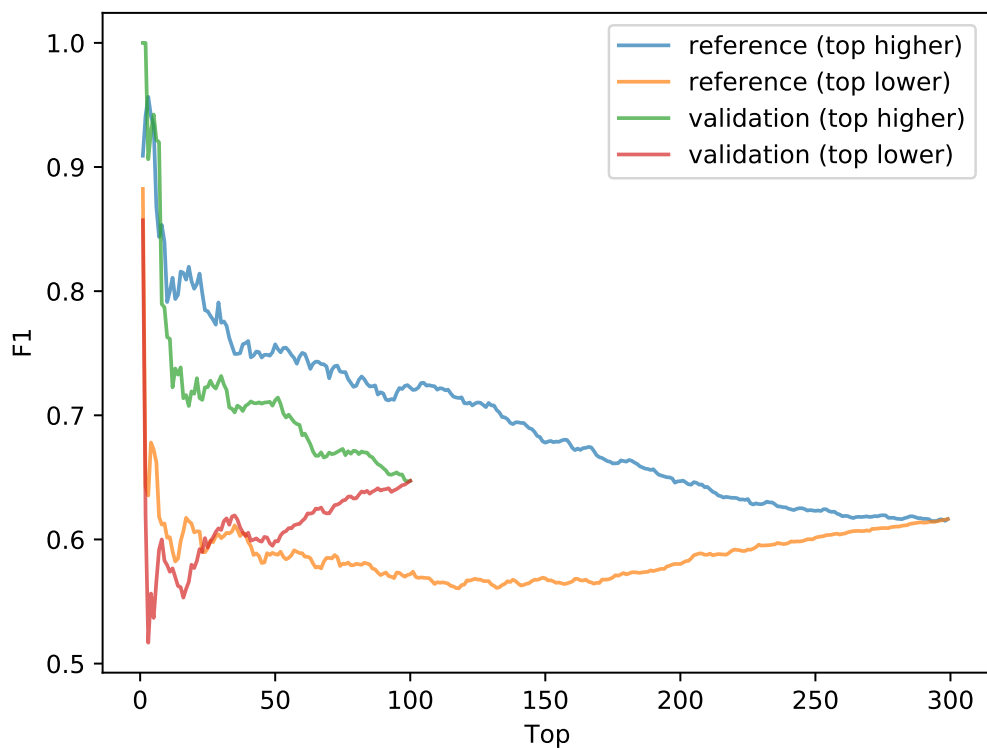


Figure 14: F_1 achieved by using only the top k sentences with higher or lower agreement respectively. The **reference (top higher)** and **reference (top lower)** lines use the sentences in the Reference Collection. Similarly, the **validation (top higher)** and **validation (top lower)** lines use the sentences in the Validation Collection. The **(top higher)** lines show the F_1 achieved by selecting only the top x sentences with higher agreement. Similarly, the **(top lower)** lines show the F_1 achieved by selecting only the top x sentences with lower agreement. As the number of sentences decreases, the collections with a higher agreement have their F_1 increased while the collections with a lower agreement have their F_1 decreased.

663 *4.2. Quality of the Ensembling Algorithms*

664 The ensembling pipeline presented in this paper was directly applied to the Blank
665 Collection of sentences.¹⁵ The ensembler was tuned in the Reference Collection and
666 validated in the Validation Collection.

667 Table 7 summarizes the performance obtained by several ensembler configurations
668 in the Reference Collection and the Validation Collection. Additionally, both the
669 baseline annotation algorithm and the submission that achieved the best score in
670 the challenge [23] (TALP-UPC [30]) are included at the start of the table. Also,
671 the oracle ensemblers are added at the end of the table for comparative purposes
672 (e.i., they provide an upper bound of the performance that can achieve a perfect
673 ensembler). Both the binary and non-binary version of the oracle ensemble achieve
674 the same approximate result.

675 The system that resulted from optimizing the combination of handcrafted strate-
676 gies achieves a slightly better performance than by only using the highest scored
677 submission (0.650 and 0.639, respectively). Table 8 summarizes the architecture of
678 the best combination found. The use of expert information through the inclusion
679 of expert weighting over uniform weighting resulted in an important improvement.
680 The optimization algorithm ran for about 500 generations of 10 individuals before
681 finishing due to no-improvement early stopping.

682 Most automatic learning systems achieved a superior performance in the Reference
683 Collection than in the Validation Collection. This was expected since weights were
684 directly optimized to maximize the performance in the training collection, i.e., the
685 Reference Collection. Additional considerations were made to measure their general-
686 ization capability, like using a “leave one out” cross-validation process. The validation

¹⁵<https://github.com/jpconsuegra/ensemble-kd2019>

System	Strategy		Ref.	Val.
Source	baseline		0.355	0.431
	TALP-UPC		0.605	0.639
Handcrafted	majority		0.426	0.486
	union		0.357	0.377
	optimization		0.619	0.650
Learning	all in one	RF	0.672	0.622
	per category	RF	0.675	0.626
	per label	SVM	0.635	0.623
		LR	0.607	0.640
Oracle	best		0.810	0.843
	all		0.834	0.873

Table 7: Summary of the performance of some relevant ensembler configurations. The baseline algorithm and the best performing submission to the challenge are included at the beginning of the table. Column **Ref** shows the F_1 score achieved in the Reference Collection. Column **Val** shows the F_1 score achieved in the Validation Collection after configuration selection. The strategy labeled as *optimization* corresponds to the ensembler configuration that resulted from exploring the space of handcrafted rules to find the optimal configuration according to the F_1 metric (shown in Table 8). RF, SVM, and LR stand for Random Forest, Support Vector Machine, and Logistic Regression, respectively.

687 revealed that indeed their performance was overestimated. The *best* and *weighting*
688 variations did not have a significant impact on the best performing learning systems,
689 though the inclusion of expert weighting produced a slightly positive improvement.

690 The performance achieved by the oracle ensemblers provides an upper bound of

Phase	Strategy	Parameters
Building	voting	binary True
	submissions	best False, top 6
Weighting	expert	
Scoring	expert	discrete True
Validating	constant-threshold	threshold 0.211

Table 8: Best ensembler configuration found in the space of handcrafted rules according to the F_1 metric. The Reference Collection was used for tuning the ensembling pipeline. Strategies and parameters are taken from Table 1.

691 the performance that can be achieved by any ensembling method that relies only on
692 the information contained across all the submissions to the eHealth-KD 2019 challenge.
693 The oracle ensembler that considers all submissions achieves a performance of 0.873
694 in the Validation Collection. This means that it is impossible to correctly annotate
695 13% of the Validation Collection based on the voting models. Relatively to this score,
696 the ensembler that resulted from the optimization process reaches approximately 74%
697 of the total capacity. This shows that there might not be any set of consistent rules
698 to decide which model votes to trust in each case.

699 4.3. Quality of the Ensembler Optimization

700 Team UH-MatCom designed and trained four deep-learning models in the 2020
701 edition of the eHealth-KD challenge [26]. Two of the four trained models were
702 submitted to the challenge. Additionally, the ensemble method described in this
703 paper was applied to produce a third run.

704 The best found ensembling pipeline for Scenario 2 and Scenario 3 was to use the

System	1-main	2-taskA	3-taskB	4-transfer
<i>top performance</i>	0.666	0.825	0.633	0.584
<i>run 1: cnet</i>	0.552	0.792	0.306	0.367
<i>run 2: deptree</i>	0.395 ^b	0.795	0.545	0.138 ^b
<i>run 3: ensemble</i>	0.557	0.795	0.005	0.373
<i>baseline</i>	0.395 ^b	0.542 ^b	0.131 ^b	0.138 ^b

Table 9: Summary of the performance (F_1) achieved by the system of team *UH-MatCom* across the three runs of the eHealth-KD 2020 challenge. The best score achieved in each scenario of the challenge is included for comparison purposes. Also, the baseline implementation provided by the challenge organizers is included. Runs that use the baseline solution are marked with “*b*”.

705 automatic learning rules, with a *per label model* configuration using a Support Vector
706 Machine (SVM) binary classifier for each different type of entity and relation. Only
707 two models were ensembled for Scenario 1 and 4 due to not having the corresponding
708 submissions of those models in these scenarios. The best found ensembling pipeline
709 was to use the automatic learning rules, with an *all in one model* configuration using
710 a single Logistic Regression classifier for all the entity types and relations.

711 Table 9 shows the results obtained by the three runs in the challenge. There was
712 an issue with the ensemble method in Scenario 3 (Task B), which resulted in an
713 ill-formed submission. The ensemble run obtained the best results out of the three
714 runs, except for the incorrectly submitted scenario. This shows that the ensembling
715 pipeline produces outputs that perform better than the systems given as input.

716 5. Discussion

717 The ensembling pipeline presented in this paper is meant to take advantage of
718 information extraction competitions, as a source of different annotated versions of the
719 same sentences, to produce an extended version of the original corpus. This method
720 responds to corpus bootstrapping, the technique in which a corpus is automatically
721 extended using programmatic instructions. The use of ensembling techniques instead
722 of using a single version has several advantages, such as: reducing the bias towards
723 a particular architecture or set of rules and providing richer information about the
724 quality of the resulting corpora. The effectiveness of the proposed method in achieving
725 these goals was evaluated in Section 4.

726 Since the ensembler presented in this paper works directly on the annotations
727 instead of working with the algorithms, it allows the integration of any group of systems
728 independently of their architecture. As long as they share the same annotation scheme
729 (which specifies the possible labels of entities and relations and the lexicographic
730 rules of the tokens), there is no need for them to share the same intermediate steps,
731 architecture, or resources. That turns scientific competitions into an opportunity for
732 automatic corpus extension since all the submitted systems share the same annotation
733 scheme.

734 Another advantage that derives from working directly on the outputs, is that
735 the ensembler can be used to combine non-programmatic sources of annotations, i.e.
736 human annotations. This advantage is especially relevant when there is no initial
737 collection of labeled data that can be used to train learning models. Unlike algorithms,
738 humans can begin annotating documents from the specification of an annotation
739 scheme. By evaluating their performance on a small trial collection, the ensembling
740 pipeline can measure the quality of their annotation as it would with computer systems.

741 Future annotations made by the same person can be intelligently combined with other
742 annotated versions of the same document to produce an improved variant. This way,
743 people from all around the world can collaborate with the construction of corpora by
744 incrementally providing their own manually annotated versions of arbitrary subsets
745 of the corpus under construction. This reduces the need for expertise or supervision
746 of the annotation task, since the ensembler can be used as a way to regularize the
747 annotation quality. Applying such a mechanism constitutes an important step into
748 the crowdsourcing corpora annotation, since more people can contribute with a little
749 supervision.

750 Compared to active learning techniques, the ensembling approach does not require
751 expert knowledge, but instead it takes advantage of having numerous sources of
752 information to build the corpus, although they may not be completely accurate. The
753 precision is regulated by the quality of the systems involved in the annotation and
754 the coverage that the creator of the corpus wishes to give to the corpus. Most of the
755 current research in active learning field involves finding the best method to choose
756 the samples to be asked to a human oracle. The ensembling pipeline applied for
757 crowdsourcing allows the annotation of arbitrary amounts of data given the softer
758 restriction on quality and the higher availability of pseudo-oracles.

759 **6. Conclusions and Future Work**

760 This paper discusses several strategies for building an ensembling pipeline oriented
761 toward NER and relation extraction tasks. Using an optimization search algorithm,
762 the space of possible ensemblers that can be built with the given strategies is explored
763 to find the optimal configuration for a concrete corpus. The *eHealth-KD v2* corpus [22]
764 was taken as a case of study, resulting in the automatic annotation of the relevant
765 key phrases and relations present across 8 000 sentences. The submissions presented

766 in the eHealth Knowledge Discovery Challenge at IberLEF 2019 were used as the
767 only source of votes (no additional voting models were trained). The experiments
768 show that the collection of ensembled sentences contains information that can be
769 used by learning algorithms to improve their performance in the original eHealth-KD
770 tasks. This fact encouraged the distribution of the collection as an auxiliary training
771 corpus in the 2020 edition of the eHealth-KD challenge. Since error correction in the
772 eHealth-KD corpus is faster than manual annotation, further progress can be made
773 now in extending it. The agreement metric computed for the sentences was shown
774 to hold information about the quality of the annotation. Additionally, the results
775 show that the ensembling pipeline produces more accurate outputs than the systems
776 it takes as input. This kind of ensembling technology can be used to produce highly
777 competitive information extraction systems.

778 Since the ensembling techniques discussed in this paper rely only on the resulting
779 annotations and not on the systems themselves, these methods can be used to assist
780 crowd annotation of NER and relation extraction related corpora. Large collections
781 of corpora can be annotated by multiple individuals with overlapping annotations. A
782 smaller collection of gold annotated sentences can be built to estimate the quality of
783 each annotation. The ensemble system will then be asked to find the best way to
784 combine the different votes to obtain a more accurate corpus. As major future work
785 we plan to study the application of machine learning techniques to deal with corpora
786 bias during the ensembling procedure to obtain a impartial version of it.

787 **Acknowledgments**

788 **Funding:** This research has been partially funded by the University of Alicante
789 and the University of Havana, the Generalitat Valenciana (*Conselleria d'Educació,*
790 *Investigació, Cultura i Esport*) and the Spanish Government through the projects

791 LIVING-LANG (RTI2018-094653-B-C22) and SIIA (PROMETEO/2018/089, PROM-
792 ETEU/2018/089). Moreover, it has been backed by the work of both COST Actions:
793 CA19134 - “Distributed Knowledge Graphs” and CA19142 - “Leading Platform for
794 European Citizens, Industries, Academia and Policymakers in Media Accessibility”.

795 **References**

- 796 [1] J. P. Chiu, E. Nichols, Named entity recognition with bidirectional lstm-cnns,
797 Transactions of the Association for Computational Linguistics 4 (2016) 357–370.
- 798 [2] L. Gligic, A. Kormilitzin, P. Goldberg, A. Nevado-Holgado, Named entity recog-
799 nition in electronic health records using transfer learning bootstrapped neural
800 networks, Neural Networks 121 (2020) 132–139.
- 801 [3] B. Agarwal, N. Mittal, Machine learning approach for sentiment analysis, in:
802 Prominent feature extraction for sentiment analysis, Springer, 2016, pp. 21–45.
- 803 [4] D. Bahdanau, K. Cho, Y. Bengio, Neural machine translation by jointly learning
804 to align and translate, arXiv preprint arXiv:1409.0473 (2014).
- 805 [5] G. Lample, M. Ballesteros, S. Subramanian, K. Kawakami, C. Dyer, Neural
806 architectures for named entity recognition, arXiv preprint arXiv:1603.01360
807 (2016).
- 808 [6] N. Limsopatham, N. Collier, Bidirectional lstm for named entity recognition in
809 twitter messages (2016).
- 810 [7] P. J. Gorinski, H. Wu, C. Grover, R. Tobin, C. Talbot, H. Whalley, C. Sudlow,
811 W. Whiteley, B. Alex, Named entity recognition for electronic health records:

- 812 a comparison of rule-based and machine learning approaches, arXiv preprint
813 arXiv:1903.03985 (2019).
- 814 [8] O. Hegazy, O. S. Soliman, M. A. Salam, A machine learning model for stock
815 market prediction, arXiv preprint arXiv:1402.7351 (2014).
- 816 [9] A. Piad-Morffis, Y. Gutiérrez, R. Muñoz, A corpus to support ehealth knowledge
817 discovery technologies, *Journal of biomedical informatics* 94 (2019) 103172.
- 818 [10] E. K. Ringger, M. Carmen, R. Haertel, K. D. Seppi, D. Lonsdale, P. McClanahan,
819 J. L. Carroll, N. Ellison, Assessing the costs of machine-assisted corpus annotation
820 through a user study., in: *LREC*, Vol. 8, 2008, pp. 3318–3324.
- 821 [11] J. Zavrel, W. Daelemans, Bootstrapping a tagged corpus through combination
822 of existing heterogeneous taggers, arXiv preprint cs/0007018 (2000).
- 823 [12] O. Ikechukwu E, O. Ebele G, A. Godwin E, et al., Bootstrapping method for
824 developing part-of-speech tagged corpus in low resource languages tagset-a focus
825 on an african igbo, arXiv preprint arXiv:1903.05225 (2019).
- 826 [13] R. Polikar, Ensemble based systems in decision making, *IEEE Circuits and*
827 *systems magazine* 6 (3) (2006) 21–45.
- 828 [14] P. Yang, Y. Hwa Yang, B. B Zhou, A. Y Zomaya, A review of ensemble methods
829 in bioinformatics, *Current Bioinformatics* 5 (4) (2010) 296–308.
- 830 [15] T. G. Dietterich, Ensemble methods in machine learning, in: *International*
831 *workshop on multiple classifier systems*, Springer, 2000, pp. 1–15.
- 832 [16] R. E. Schapire, The strength of weak learnability, *Machine learning* 5 (2) (1990)
833 197–227.

- 834 [17] L. Breiman, Bagging predictors, *Machine learning* 24 (2) (1996) 123–140.
- 835 [18] L. Rokach, Ensemble-based classifiers, *Artificial Intelligence Review* 33 (1-2)
836 (2010) 1–39.
- 837 [19] D. Opitz, R. Maclin, Popular ensemble methods: An empirical study, *Journal of*
838 *artificial intelligence research* 11 (1999) 169–198.
- 839 [20] R. Speck, A.-C. N. Ngomo, Ensemble learning for named entity recognition, in:
840 *International semantic web conference*, Springer, 2014, pp. 519–534.
- 841 [21] A. Piad-Morffis, Y. Gutiérrez, Y. Almeida-Cruz, R. Muñoz, [dataset] *eHealth-KD*
842 *v2* (Mar. 2020). doi:10.5281/zenodo.3696792.
843 URL <https://doi.org/10.5281/zenodo.3696792>
- 844 [22] A. Piad-Morffis, Y. Gutiérrez, Y. Almeida-Cruz, R. Muñoz, *A compu-*
845 *tational ecosystem to support ehealth knowledge discovery technologies*
846 *in spanish*, *Journal of Biomedical Informatics* (2020) 103517doi:<https://doi.org/10.1016/j.jbi.2020.103517>.
847 [//doi.org/10.1016/j.jbi.2020.103517](https://doi.org/10.1016/j.jbi.2020.103517).
848 URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S1532046420301453)
849 [S1532046420301453](http://www.sciencedirect.com/science/article/pii/S1532046420301453)
- 850 [23] A. Piad-Morffis, Y. Gutiérrez, J. P. Consuegra-Ayala, S. Estevez-Velarde,
851 Y. Almeida-Cruz, R. Munoz, A. Montoyo, Overview of the ehealth knowledge
852 discovery challenge at iberlef 2019, in: *Proceedings of the Iberian Languages*
853 *Evaluation Forum (IberLEF 2019)*. CEUR Workshop Proceedings, CEUR-WS.
854 org, 2019.
- 855 [24] J. P. Consuegra, [jpconsuegra/ehealthkd-2019-ensembled-corpus](https://github.com/jpconsuegra/ehealthkd-2019-ensembled-corpus): Version 20.07

- 856 (Jul. 2020). doi:10.5281/zenodo.3926746.
857 URL <https://doi.org/10.5281/zenodo.3926746>
- 858 [25] A. Piad-Morffis, Y. Gutiérrez, S. Estevez-Velarde, Y. Almeida-Cruz, R. Muñoz,
859 A. Montoyo, Overview of the eHealth Knowledge Discovery Challenge at IberLEF
860 2020, in: Proceedings of the Iberian Languages Evaluation Forum (IberLEF
861 2020), 2020.
- 862 [26] J. P. Consuegra-Ayala, M. Palomar, UH-MatCom at eHealth-KD Challenge 2020,
863 in: Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2020),
864 2020.
- 865 [27] A. Piad-Morffis, J. P. Consuegra-Ayala, Y. Gutiérrez, S. Estevez-Velarde, [ehealth-](#)
866 [kd 2020](#) (Apr. 2020). doi:10.5281/zenodo.3870509.
867 URL <https://doi.org/10.5281/zenodo.3870509>
- 868 [28] L. Ratinov, D. Roth, Design challenges and misconceptions in named entity
869 recognition, in: Proceedings of the Thirteenth Conference on Computational
870 Natural Language Learning (CoNLL-2009), 2009, pp. 147–155.
- 871 [29] J. Lafferty, A. McCallum, F. C. Pereira, Conditional random fields: Probabilistic
872 models for segmenting and labeling sequence data (2001).
- 873 [30] S. Medina Herrera, J. Turmo Borrás, Talp-upc at ehealth-kd challenge 2019:
874 A joint model with contextual embeddings for clinical information extraction,
875 in: Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2019):
876 co-located with 35th Conference of the Spanish Society for Natural Language
877 Processing (SEPLN 2019): Bilbao, Spain, September 24th, 2019, CEUR-WS.
878 org, 2019, pp. 78–84.