

Document downloaded from:

<http://hdl.handle.net/10251/171314>

This paper must be cited as:

Hernandez-Farias, D.I.; Prati, R.; Herrera, F.; Rosso, P. (2020). Irony Detection in Twitter with Imbalanced Class Distributions. *Journal of Intelligent & Fuzzy Systems*. 39(2):2147-2163. <https://doi.org/10.3233/JIFS-179880>



The final publication is available at

<https://doi.org/10.3233/JIFS-179880>

Copyright IOS Press

Additional Information

# Irony detection in Twitter with imbalanced class distributions

Delia Irazú Hernández Farías<sup>a,\*</sup>, Ronaldo Prati<sup>b</sup>, Francisco Herrera<sup>c</sup> and Paolo Rosso<sup>d</sup>

<sup>a</sup>*División de Ciencias e Ingenierías Campus León, Universidad de Guanajuato, Mexico*

<sup>b</sup>*Universidade Federal do ABC, Brazil*

<sup>c</sup>*Department of Computer Science and Artificial Intelligence, University of Granada, Spain*

<sup>d</sup>*Universitat Politècnica de València, Spain*

**Abstract.** Irony detection is a not trivial problem and can help to improve natural language processing tasks as sentiment analysis. When dealing with social media data in real scenarios, an important issue to address is data skew, i.e. the imbalance between available ironic and non-ironic samples available. In this work, the main objective is to address irony detection in Twitter considering various degrees of imbalanced distribution between classes. We rely on the emotIDM irony detection model. We evaluated it against both benchmark corpora and skewed Twitter datasets collected to simulate a realistic distribution of ironic tweets. We carry out a set of classification experiments aimed to determine the impact of class imbalance on detecting irony, and we evaluate the performance of irony detection when different scenarios are considered. We experiment with a set of classifiers applying class imbalance techniques to compensate class distribution. Our results indicate that by using such techniques, it is possible to improve the performance of irony detection in imbalanced class scenarios.

**Keywords:** Irony detection, class imbalance, imbalanced learning

## 1. Introduction

Users of social media platforms tend to formulate points of view, opinions, and judgments concerning almost everything surrounding them: from a given event up to a personal experience. Social media allow the users to employ language in its literal sense but sometimes figurative language devices are also exploited. Among them, there is one that serves to express opinions in a witty (and often funny) way: *irony*.

Irony serves to express an evaluative judgment or attitude towards a particular target [2]. It allows us to convey subjective ideas by using the non-literal

meaning of the words. Several theories have been proposed attempting to describe what irony is. Perhaps the most common one is that from the Gricean tradition [14], where irony is defined as a trope where the speaker intends to communicate the opposite meaning of what is literally said. Other authors consider that an ironic utterance serves to reveal the speaker's position (approval or disapproval) on the result of something [25, 43].

Besides the different theories defining irony, such a language device is also related to another concept: sarcasm. Both terms are perceived as synonyms due to the subtle distinction between them. When irony involves stressed negative evaluation towards a particular target with the intention of a given offense, it is considered as sarcasm [4, 30]. Under a computational linguistic perspective, irony and sarcasm are considered either as synonyms or being irony an umbrella term covering sarcasm.

In social media, people use irony, having (most of the time) only an intuitive definition of this concept. Consequently, the ironic content in these platforms reflects what people consider as this kind of figurative language device. As can be noticed in the following tweets<sup>1</sup>, irony can be used with different purposes: to express an evaluation indirectly (example (i)) or to reveal a failed expectation (example (ii)).

- (i) I seriously loveeeee how much you care  
(ii) My train got cancelled.. Good way to start the day! -.- #Västtrafik

Interest in detecting the presence of irony in social media has grown significantly in the past years. Understanding the real meaning of a given message is an ongoing task for computational linguistics; therefore, such an intriguing figurative language device represents a big challenge.

In particular, Twitter data have become popular for irony detection [21]. Twitter represents an interesting source of information regarding how people perceive events, products, and so on. It provides a huge amount of user-generated data (easily accessible via Twitter API<sup>2</sup>) that allows capturing a wide variety of real uses of irony in this kind of short texts. Several approaches have been proposed to deal with irony detection relying on different perspectives. The authors in [6, 23, 40] addressed irony detection relying mainly on textual-based features. Others [5, 24, 44] took advantage of information regarding the context surrounding a given comment to determine whether or not an ironic meaning is intended. Exploiting the affective property of irony, in [17] the authors proposed an approach considering mainly such kind of information. Moreover, novel techniques such as word-embeddings and deep learning models have also been exploited [22, 34,36].

Despite data skew has been recognized as a critical issue for irony detection [21], related work addressing this task as a class imbalance problem is scarce. The Multi-Strategy Ensemble Learning Approach (MSELA) was proposed by Liu et al. [28] to deal with irony detection in imbalanced class datasets. The authors experimented with ironic comments written in English and Chinese. The MSELA combines sample-ensemble, classifier-ensemble, and weighted voting strategies together with a set of different features for each language. Punctuation marks, n-grams, and POS tags were used as features for English,

whereas extreme positive and negative nouns, adjectives, adverbs of degree and proverbs were exploited for Chinese. Results on different settings exploiting MSELA were reported achieving in overall 0.8 in AUC (Area Under the ROC curve) terms.

A corpus of manually annotated Twitter conversation was used by Abercrombie and Hovy [1] for experimenting with irony detection on balanced and imbalanced class scenarios. Furthermore, the authors compare the performance of recognizing the irony of both human and machine learning algorithms. A logistic regression classifier with features such as n-grams and POS tags was employed. The performance of the proposed approach suffers from significant drops on the imbalanced class data in both F-measure and AUC terms when compared with those from the balanced one. Cervone et al. [10] experimented with ironic tweets written in Italian by applying balancing techniques to address the data skew. They exploited different sets of features combined with random oversampling, undersampling, and cost-sensitive learning. The best performance was obtained by the last one when it was used together with bag-of-words representation.

The progress so far achieved in irony detection has been focused on the development of models able to automatically capture potential cues for identifying such kind of figurative language device. However, even when the data skew has been recognized as an inherent factor related to the presence of ironic content on Twitter, the majority of the related work fails to have regard to the role of imbalanced class degree.

In this paper, we address irony detection from a perspective of imbalanced distributions aim at evaluating the impact of applying different preprocessing techniques for detecting irony in Twitter. As far as we know, this is the first work where irony detection is addressed by considering many factors related to the imbalanced learning problem. It is important to highlight that we are not proposing a novel technique for dealing with imbalanced class data, instead, we are experimenting with the use of existing techniques for evaluating irony detection in an imbalanced class scenario. We are considering several factors related to the class imbalance problem described by [20]. Another important point arises from the fact that we are not introducing a new approach to detect irony rather, we are using the model described in [17]. Furthermore, we carried out an extensive experimental study with a large benchmark of irony detection corpora that covers several aspects ranging from developing criteria to imbalanced class degree.

Summarizing, our main contributions are the following:

1. We exploited the irony detection model described in [17] in order to carry out a set of classification experiments aimed to determine the impact of class imbalance for detecting irony.
2. We experimented with several treatment techniques in order to assess its impact on the performance of irony detection.
3. We developed a new set of corpora (denoted as *TwImbData*) retrieved considering criteria for simulating a realistic scenario. This dataset could serve as a starting point for expanding the research on irony detection considering an imbalanced class scenario.

*Organization.* The paper is organized as follows. Section 2 introduces irony detection from a class imbalanced problem perspective; besides, the irony detection model and the corpora we used for experimental purposes are also presented. In Section 3 we describe the experimental setting and the obtained results of addressing irony detection as a class imbalance problem. Section 4 summarizes the main findings of our study. Finally, in Section 5 we draw some conclusions and some directions for future work.

## 2. Irony Detection as an imbalanced class classification problem

Every day millions of tweets are posted on Twitter<sup>3</sup>. Even when the use of irony in this social media platform is quite common, the difference between the amount of non-ironic, i.e., tweets expressing any other kind of intention, and ironic tweets is enormous. Therefore, an important issue to address in irony detection is data skew, i.e., the imbalance between ironic and non-ironic samples available, as it reflects the realistic distribution of the use of irony in Twitter [1, 27, 39]. This problem also happens in many other real-world problems such as biology, medicine, economy, etc. The role of data skew for detecting the presence of irony in social media has been recognized as an important challenge [1, 21] that needs to be considered when designing irony detection models.

Furthermore, according to [3], a class imbalance distribution problem could occur in two situations: (i)

when class imbalance occurs naturally in the problem in hand, and (ii) when the data is not imbalanced by definition, instead is very expensive to acquire such data for minority class due to factors such as cost, effort, etc. The irony detection problem fits with both situations. First, there is a big difference between the ironic and non-ironic tweets published in a given time frame. Secondly, retrieving potentially ironic data from a given data source is not a trivial task. Two different methodologies for acquiring data for irony detection have been recognized in [18]: self-labeling and crowdsourcing. The former one involves the use of certain labels such as hashtags that are added by the authors of the texts, while in the second one a manual annotation procedure needs to be performed; then, the task becomes more complex involving the inherent subjectivity of irony not only from the author of the text in hand but also from the human annotator.

Generally speaking, irony detection has been addressed as a binary classification task. The main aim is to distinguish ironic from non-ironic texts. In a nutshell, when an irony detection approach is proposed, the principal goals are: (i) to propose a set of relevant features helping to capture the ironic intention in a given text and, (ii) to assess the performance of the model usually in an in-house dataset collected by the authors. Most of these approaches do not consider other related aspects such as the impact of the inherent imbalanced nature of the presence of irony in social media platforms.

In the next section, we introduce in detail the irony detection model we exploited for detecting ironic tweets in imbalanced class scenarios.

### 2.1. *emoIDM*: an irony detection model based on affective information

According to several theorists, affect plays an important role in the use of irony [2, 47]. Therefore, considering the presence of affective content involved in ironic texts represents an interesting starting point. Attempting to take advantage of such kind of information, we rely on *emoIDM* proposed in [17].

*emoIDM* addresses irony detection as a classification task by considering different facets of affective content as well as structural markers. To represent a tweet, *emoIDM* uses three different groups of features:

1. *Structural*. It includes punctuation marks, length of words, part-of-speech labels, Twitter Marks (i.e., hashtags, mentions, etc.), among others.

- 245 2. *Sentiment*. In order to capture sentiment information, *emotIDM* takes advantage of a wide  
 246 range of English lexical resources such as:  
 247 AFINN [33], Hu&Liu [19], among others.  
 248
- 249 3. *Emotions*. Attempting to cover as much  
 250 information related to emotions as possible,  
 251 *emotIDM* considers features regarding the  
 252 main theories in the nature of emotions: Cate-  
 253 gorical and Dimensional models of emotions.  
 254 The Categorical model suggests the existence  
 255 of basic emotions such as anger, fear, joy, dis-  
 256 gust, etc., that in *emotIDM* are considered from:  
 257 EmoLex [31], EmoSenticNet [37], and LIWC  
 258 [35]. In the Dimensional model, an emotion can  
 259 be defined according to its position in a space  
 260 of independent dimensions. *emotIDM* includes  
 261 the dimensions defined in: ANEW [8], Diction-  
 262 ary of Affect in Language [46], and SenticNet  
 263 [9].

264 Exploiting affective information for detecting  
 265 irony also allows to capture this kind of informa-  
 266 tion, disregarding domain. Besides, in line with most  
 267 of current approaches in computational linguistics,  
 268 irony here is considered as an umbrella term that also  
 269 covers sarcasm. Tackling differences between these  
 270 devices in social media is a further challenge for fig-  
 271 urative language processing [42, 45], which is very  
 272 interesting but beyond the scope of this work.

273 Most of the time, when an irony detection model  
 274 is proposed, it is evaluated over an in-house dataset  
 275 retrieved by its authors. Instead, the performance of  
 276 *emotIDM* was assessed by using a set of corpora in  
 277 the state of the art. The obtained results outperformed  
 278 those in the related work and validated the importance  
 279 of affect-related information for detecting ironic con-  
 280 tent in tweets.

## 281 2.2. Irony detection corpora

282 In a similar fashion than in other natural language  
 283 processing tasks, collecting user-generated data con-  
 284 taining ironic instances is not a simple task. As  
 285 mentioned before, two main approaches have been  
 286 adopted for collecting Twitter data:

- 287 i By taking advantage of hashtags (such as  
 288 “#irony” and “#sarcasm”) that allow users to  
 289 explicitly marking their tweets as ironic. The  
 290 readability of using hashtags as golden labels has  
 291 been experimentally confirmed [26].  
 292 ii By exploiting crowdsourcing techniques to deter-  
 293 mine whether a tweet is ironic or not.

294 Interest in investigating the use of irony in Twitter  
 295 has led into having a wide set of available corpora for  
 296 addressing irony detection. Nevertheless, there are  
 297 not specific corpora developed considering imbal-  
 298 anced class scenarios, i.e., a dataset which keeping  
 299 the inherent imbalanced class ratio of this problem in  
 300 a real scenario is considered. We experimented with  
 301 two different groups of corpora: (a) *Benchmark cor-  
 302 pora*, and (b) *Imbalanced Class Twitter data for Irony  
 303 Detection*. Next, we describe both groups of corpora  
 304 as well as its main characteristics.

### 305 *Benchmark corpora*

306 As mentioned before, there are several Twit-  
 307 ter corpora developed for evaluating different irony  
 308 detection approaches. We took advantage of the five  
 309 corpora described below:

- 310 – *TwReyes2013*. Reyes et al. [40] collected a set of  
 311 tweets by taking advantage of specific hashtags.  
 312 They selected three hashtags for collecting non  
 313 ironic tweets: #education, #humor, and #politics.  
 314 Concerning to the ironic instances, they relied on  
 315 the use of the hashtag #irony by Twitter users.
- 316 – *TwRiloff2013*. Riloff et al. [41] created a  
 317 Twitter corpus of 3,200 tweets following a  
 318 hybrid approach involving the presence of spe-  
 319 cific markers as well as crowdsourcing. They  
 320 retrieved tweets containing sarcastic hashtags  
 321 (such as #sarcasm and #sarcastic) and also some  
 322 regular tweets. Then, they asked three annotators  
 323 to manually annotate the presence of sarcastic  
 324 content in the tweets after removing the afore-  
 325 mentioned hashtags (if any).
- 326 – *TwBarbieri2014*. Barbieri et al. [6] adopted a  
 327 similar methodology to the one of [40]. The  
 328 non ironic tweets are composed by those equiv-  
 329 alents in the *TwReyes2013* together with 10,000  
 330 tweets collected by exploiting the #newspaper  
 331 hashtag. Regarding the ironic tweets, the authors  
 332 took advantage of two hashtags: #irony and  
 333 #sarcasm<sup>4</sup>.
- 334 – *TwPtáček2014*. Ptáček et al. [39] introduced  
 335 two sarcastic datasets: in Czech<sup>5</sup> and English.  
 336 For collecting the sarcastic tweets the authors  
 337 used the hashtag #sarcasm, while the non sar-  
 338 castic tweets were collected using only the

<sup>4</sup>In the rest of the paper, we will use *TwIronyBarbieri2014* and *TwSarcasmBarbieri2014* to refer which set of tweets are used as ironic tweets those with #irony or #sarcasm, respectively.

<sup>5</sup>More details about this dataset can be found in [39].

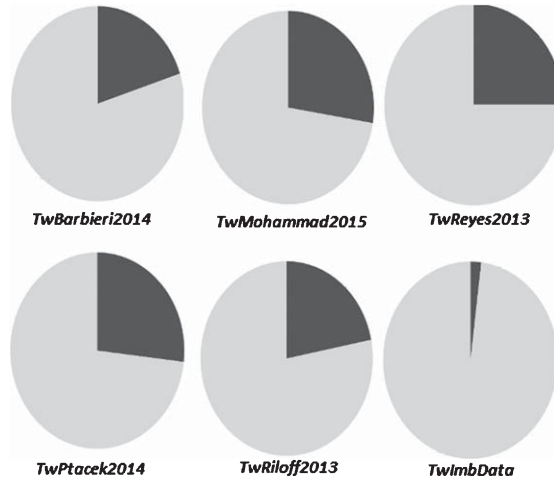


Fig. 1. “Ironic” and “non-ironic” tweets distribution in the corpora.

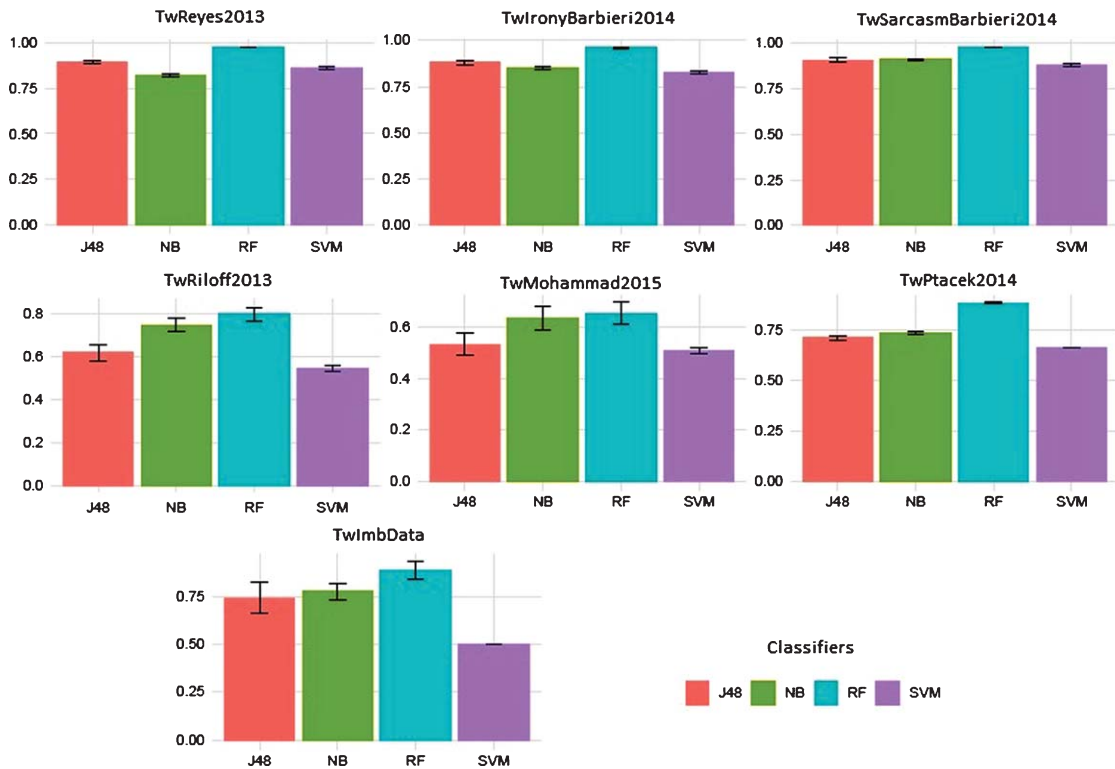


Fig. 2. Obtained results in AUC terms using the original distribution of the corpora.

339  
340  
341

language (English) as a filter parameter. The *TwPtacek2014* comprises two different distribution scenarios: balanced, and imbalanced<sup>6</sup>.

<sup>6</sup>In this paper, we used a subset of the tweets in the imbalanced class distribution because of the perishability of Twitter data.

– *TwMohammad2015*. Mohammad et al. [32] collected a set of tweets related to the 2012 US presidential elections<sup>7</sup>. They defined a multi-

342  
343  
344

<sup>7</sup>Some hashtags such as #election2012, #election, #president2012, among others were used for retrieving data from Twitter.

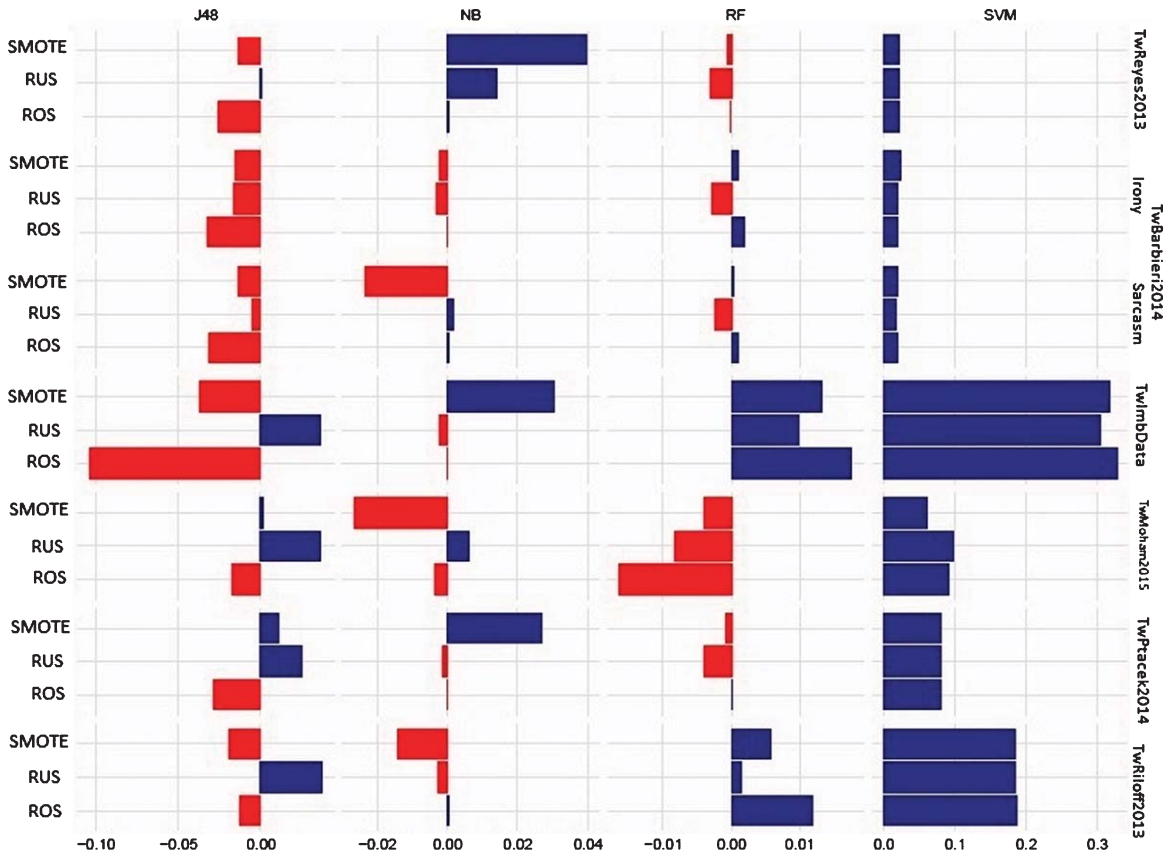


Fig. 3. Differences in terms of AUC with respect to the results on the ORIGINAL distribution after applying treatment techniques.

layer annotation schema concerning different aspects such as sentiment, emotions, purpose, and style. The last one includes sarcasm, hyperbole, understatement, and simple statement as labels.

To sum up, the *TwReyes2013*, *TwBarbieri2014*, and *TwPtáček2014* were retrieved by relying on the presence of specific labels used by the users to point out an ironic (or sarcastic) intention. Instead, *TwRiloff2013* and *TwMohammad2015* involve manual annotation of ironic tweets by exploiting crowdsourcing techniques. Regarding to *TwReyes2013* and *TwBarbieri2014*, we have merged all “non-ironic” samples into a unique class.

### ***New Imbalanced Twitter Corpora for Irony Detection***

With the aim to simulate a “realistic scenario”, i.e. a dataset that resembles a hypothetical proportion

of “ironic” tweets with respect to the “non-ironic” ones, we retrieved data from Twitter by exploiting the Streaming API. Many factors are influencing the number of tweets posted in a day. Therefore, providing a fixed or approximate quantity of ironic tweets posted in a day is not possible. We collected a sample of the tweets posted in a day (from 8th up to 18th November 2016) applying a two-step filtering criteria:

1. The tweets must contain at least one of the following hashtags: “#irony” and “#sarcasm”.
2. The tweets must be written in English.

The ironic instances were retrieved by following both criteria. Instead, the “non-ironic” instances are those tweets collected only with the second criterion. A total of eleven datasets were created, the “ironic” instances are those collected with the first criterion while in the case of the “non-ironic” we randomly selected a subset of tweets according to a

383 fixed imbalance ratio between-class of 1 : 50 (i.e., for  
 384 each “ironic” tweet, there are 50 “non-ironic” tweets).  
 385 Such datasets were grouped into a single one denoted  
 386 as *TwImbData*, where each subset has the same imbalance  
 387 ratio.

388 To sum up, Fig. 1 shows the distribution of  
 389 ironic (in black color) and non-ironic (in gray color)  
 390 tweets in the set of corpora described in this section.  
 391 As can be noticed, the distribution among  
 392 classes in *TwImbData* is very different from the other  
 393 datasets.

### 394 3. Addressing irony detection with 395 imbalanced data

396 In supervised classification, the prediction of rare  
 397 events is known as the class imbalance problem [12,  
 398 38]. Class imbalance may imply great challenges for  
 399 machine learning algorithms. Most of them tend to  
 400 misclassify the minority instances more often than  
 401 the majority instances on imbalanced class datasets.  
 402 Aimed to determine the impact of class imbalance

403 for detecting irony, we performed an experimental  
 404 setting considering several aspects.

405 We carried out a set of experiments to evalu-  
 406 ate the performance of *emotIDM* under different  
 407 degrees of class imbalance by applying differ-  
 408 ent methods for compensating class distribution.  
 409 To deal with the class imbalance, many solu-  
 410 tions have been proposed in the past few years  
 411 [15]. These solutions can be broadly categorized  
 412 into two groups: (i) *data level approaches* and (ii)  
 413 *algorithm level approaches*.

414 *Data level approaches* work in the preprocessing  
 415 phase. They are independent of the learning algo-  
 416 rithm, and in general, aim to re-balance the data  
 417 distribution by discarding (undersampling) major-  
 418 ity or replicating (oversampling) minority instances.  
 419 Simple approaches to do this include random under  
 420 sampling (hereafter RUS) and random oversampling  
 421 (denoted as ROS) [7]. There are some disadvan-  
 422 tages related to the use of these techniques, for  
 423 example, with RUS, there is a possibility of dis-  
 424 carding useful data for the learning process. On the

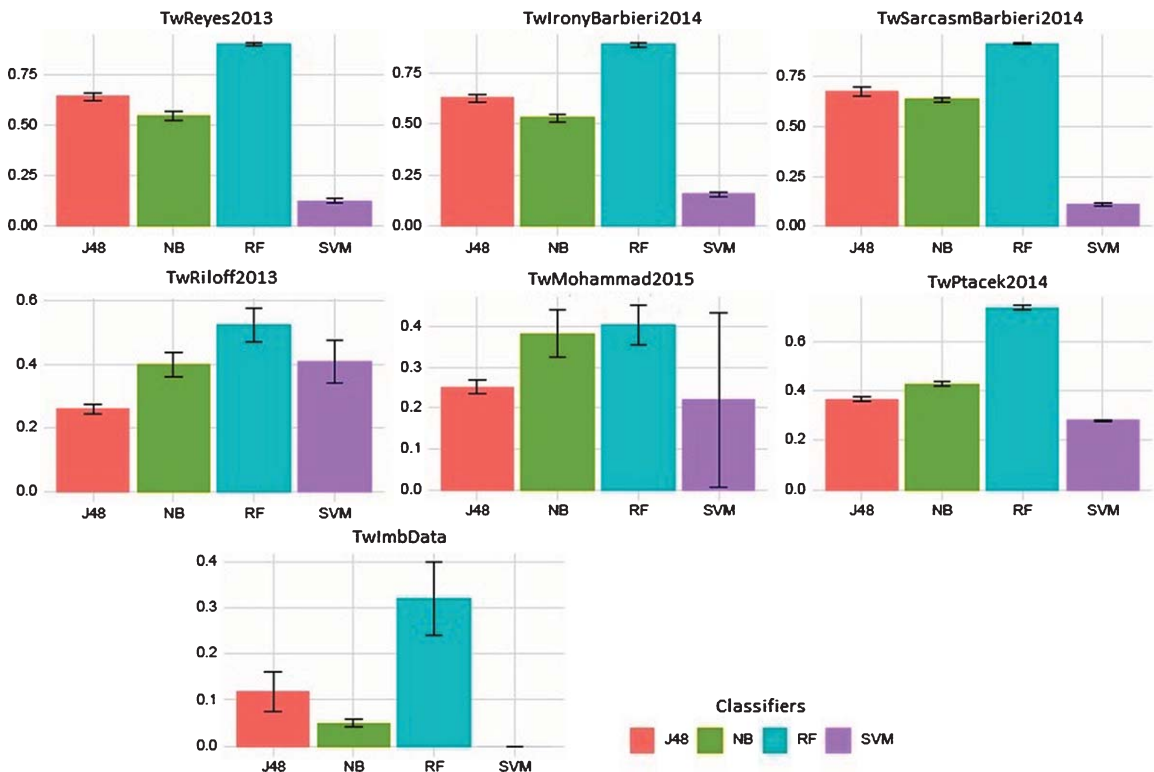


Fig. 4. Obtained results in AUPR terms using the original distribution of the corpora.



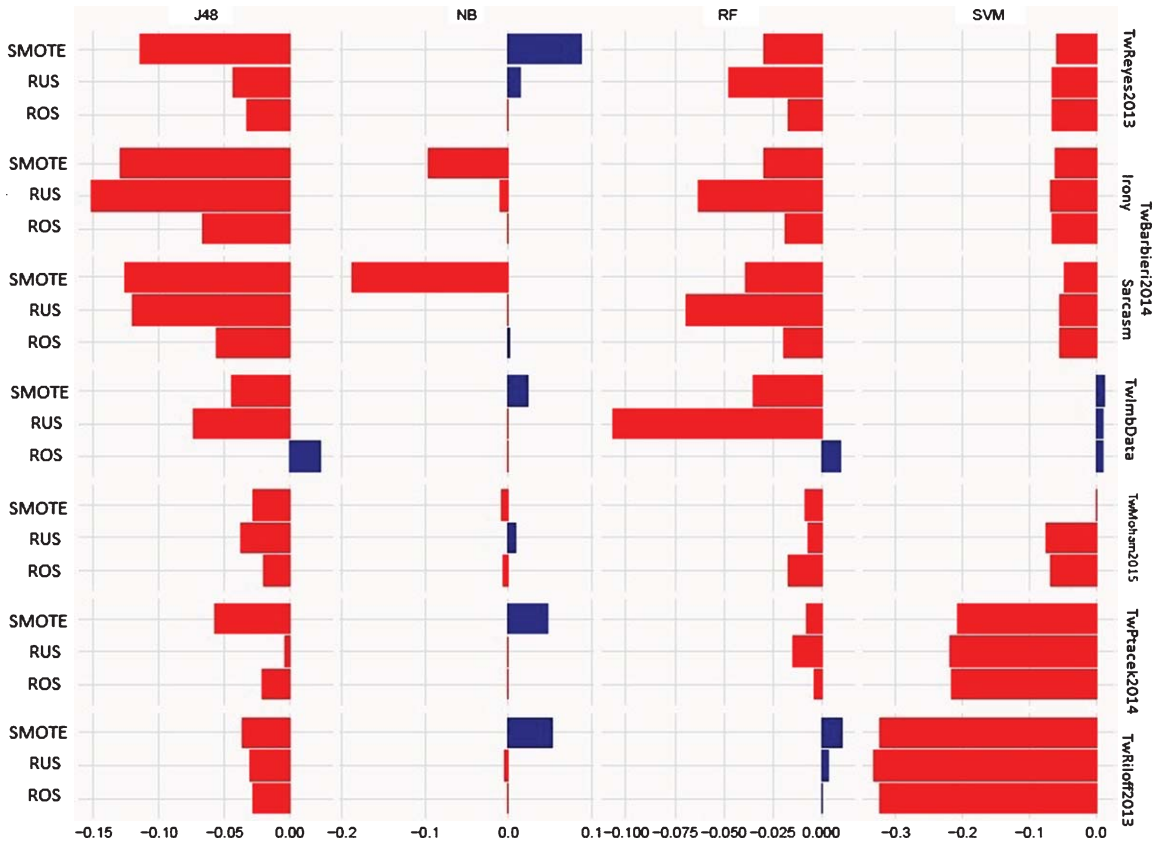


Fig. 5. Differences in terms of AUPR with respect to the results on the ORIGINAL distribution after applying treatment techniques.

425 other hand, with ROS the probability of provoking  
 426 overfitting increases. An approach that syntheti-  
 427 cally generates instances from the minority class  
 428 is the Synthetic Minority Oversampling Technique,  
 429 denoted as SMOTE [11]. The main idea of SMOTE  
 430 is to create new instances of the minority class by  
 431 interpolating them in order to oversample the train-  
 432 ing set. Apart from that, *algorithm level approaches*  
 433 involve the adaptation of learning algorithms to deal  
 434 with class imbalance. These modifications gener-  
 435 ally involve the adjustment of some optimization  
 436 criteria to trade-off frequent and infrequent classes  
 437 differently.

438 We addressed the classification between "ironic"  
 439 and "non-ironic" tweets by exploiting the Weka<sup>8</sup>  
 440 implementation of the following machine learning  
 441 classifiers (the default parameters were used for  
 442 experimental purposes): Naive Bayes (NB), Deci-  
 443 sion Tree (J48), Support Vector Machine (SVM), and  
 444 Random Forest (RF). We ran the experiments using

<sup>8</sup><http://www.cs.waikato.ac.nz/ml/index.html>

445 five-fold cross validation within each dataset from  
 446 the corpora. The experiments were paired, that is,  
 447 the same training and test partitions were used for all  
 448 learning algorithms.

449 In order to compensate for different class imbal-  
 450 ance distributions in irony detection, we applied three  
 451 class imbalance treatment techniques, namely ROS,  
 452 RUS, and SMOTE. The aforementioned data level  
 453 techniques were applied in order to achieve a bal-  
 454 anced (50% of instances in each class) proportion in  
 455 the training set. The class imbalance treatment meth-  
 456 ods were applied to the training set, and the test set  
 457 was left untouched. For the sake of comparison, we  
 458 also used the original distribution (denoted as ORIG-  
 459 INAL) as presented in each of the corpora described  
 460 in Section 2.2.

461 We are interested in assessing the performance of  
 462 irony detection when imbalance treatment techniques  
 463 are used in order to compensate for the differences  
 464 in terms of instances per class. As evaluation met-  
 465 rics we considered five different namely: Area Under  
 466 the Curve (AUC), Area Under the Precision-Recall

467 Curve (AUPR), Balanced Accuracy (BAC), Predic- 491  
 468 tive Positive (PPOS), and F1 score. Being the last one 492  
 469 the most common used for assessing the performance 493  
 470 of irony detection in Twitter. 494

471 In the following paragraphs, the obtained results 495  
 472 of applying the aforementioned experimental setting 496  
 473 are described. For each evaluation metric, we present 497  
 474 two figures with the obtained results. The first one 498  
 475 reflects the outcomes over the ORIGINAL distri- 499  
 476 bution. After applying the treatment techniques, we 500  
 477 calculate the difference between the obtained result 501  
 478 over the ORIGINAL distribution and the correspond- 502  
 479 ing performance when using a given preprocessing 503  
 480 method. Therefore, when this difference is posi- 504  
 481 tive (i.e., there is an improvement of the results), 505  
 482 it is represented as a bar towards the right side. 506  
 483 On the contrary, when the difference is negative 507  
 484 (i.e., the obtained result over the original distri- 508  
 485 bution decreased), it is represented as a bar towards the 509  
 486 left side. 510

487 Each of the *Benchmark* corpus is presented indi- 511  
 488 vidually, while in the case of the *TwImbData* we 512  
 489 present the average result of considering each dataset 513  
 490 individually. All the experiments were performed in

each of the datasets composing *TwImbData*, however 491  
 for the sake of the readability, we decided to group 492  
 the obtained results since those corpora share similar 493  
 proprieties. 494

### Area Under the Curve

Figure 2 shows the obtained results over the ORIGINAL distribution considering AUC as evaluation metric. In all corpora, the highest results were obtained using RF as the classifier. SVM emerged as the classifier with the lowest performance in the ORIGINAL distribution.

As it is shown in Fig. 3, when the treatment techniques were applied together with SVM in all corpora, there is a positive impact on the results with respect to the performance of the ORIGINAL distribution. On the other hand, there is a negative impact of using J48 with treatment techniques, except with RUS when it is used for experimental purposes on most of the corpora. Regarding the use of NB, there are some cases where using SMOTE allows improving its performance against the ORIGINAL distribution. The overall performances in terms of AUC of the imbalance treatment techniques are lower in those datasets

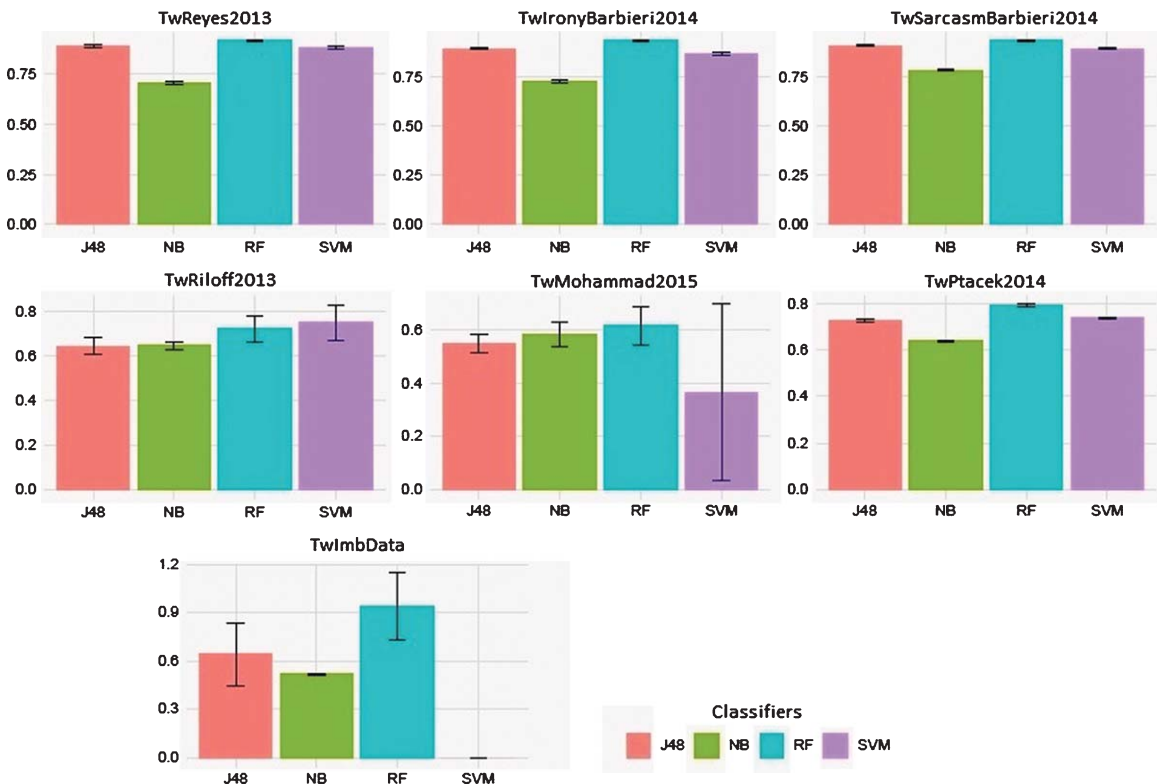


Fig. 6. Obtained results in Balanced Accuracy terms using the original distribution of the corpora.

514 where crowdsourcing was involved for developing  
 515 corpora, in line with the findings of [17]. Generally  
 516 speaking, the performance of the model in terms of  
 517 AUC across the corpora reveals an improvement in  
 518 the performance in most cases when treatment tech-  
 519 niques are applied. The lowest rates were achieved  
 520 in *TwMohammad2015* while the best ones were in  
 521 *TwBarbieri2014* and *TwReyes2013*.

522 As can be observed, the most noticeable differ-  
 523 ences are in corpora with a higher imbalanced class  
 524 degree rate. In *TwImbData* the increase is around 0.3  
 525 for all the treatment techniques. In the case of J48, in  
 526 most of the experiments, there is a negative impact  
 527 in terms of AUC. Applying treatment techniques  
 528 together with RF helps to enhance the performance  
 529 of the classifiers in *TwImbData* and *TwRiloff2013*.

530 **Area Under the Precision-Recall Curve**

531 In Fig. 4 we present the outcomes of the exper-  
 532 imental setting when AUPR was considered as

533 evaluation metric. AUPR is considered as a use-  
 534 ful measure of success of prediction when the  
 535 classes are very imbalanced, as this case. The best  
 536 performance in terms of AUPR was achieved by the  
 537 RF; while the SVM has the lowest rates.

538 In most of the cases, there is a drawback in the  
 539 performance in terms of AUPR of the classifiers  
 540 when treatment techniques were applied (as shown in  
 541 Fig. 5). Considering those experiments where there  
 542 is an improvement, it can be observed that it was  
 543 achieved by either SMOTE or ROS. In terms of  
 544 AUPR, when SVM was used the obtained results  
 545 over the benchmark corpora were not improved by  
 546 applying treatment techniques. This is not the case of  
 547 *TwImbData*, where using all the preprocessing tech-  
 548 niques there is a slight improvement with respect  
 549 to the ORIGINAL distribution. The most signifi-  
 550 cant improvement considering AUPR was obtained  
 551 when NB is used together with SMOTE in the  
 552 *TwReyes2013*.

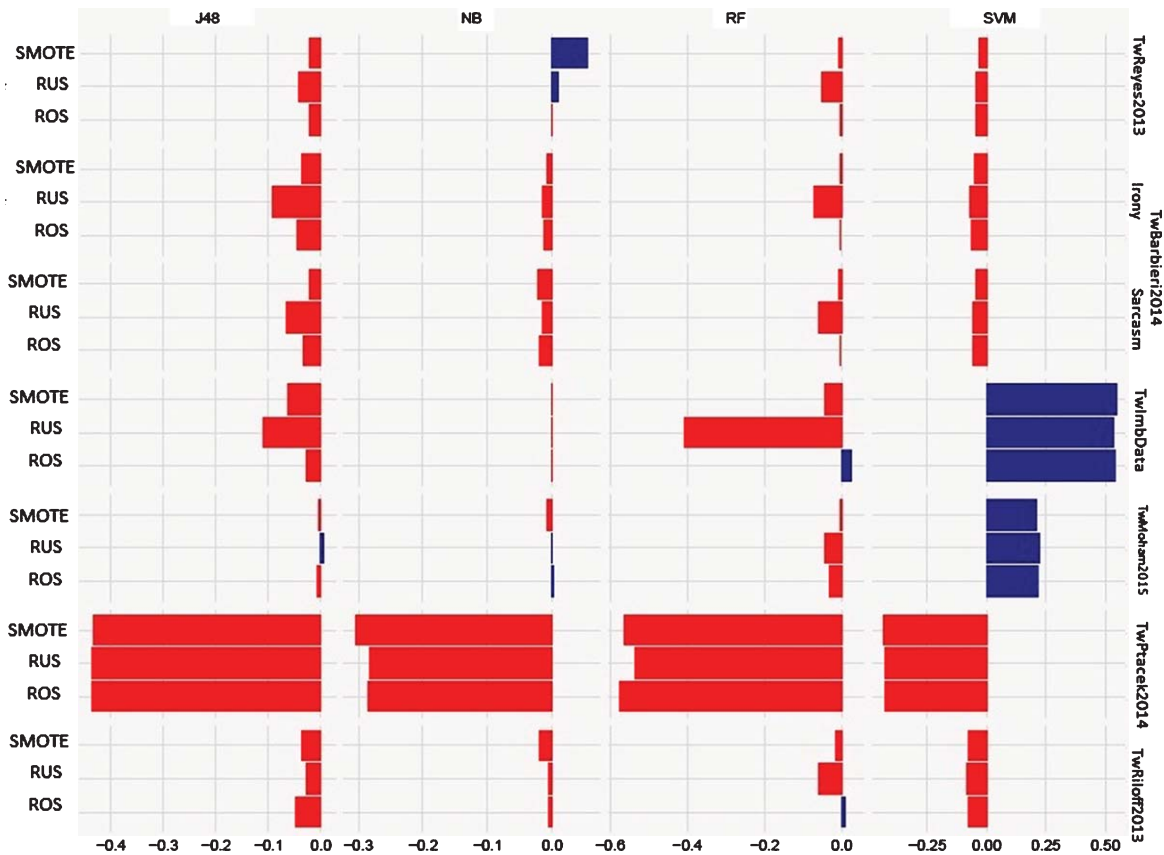


Fig. 7. Differences in terms of Balanced Accuracy with respect to the results of the ORIGINAL distribution after applying treatment techniques.

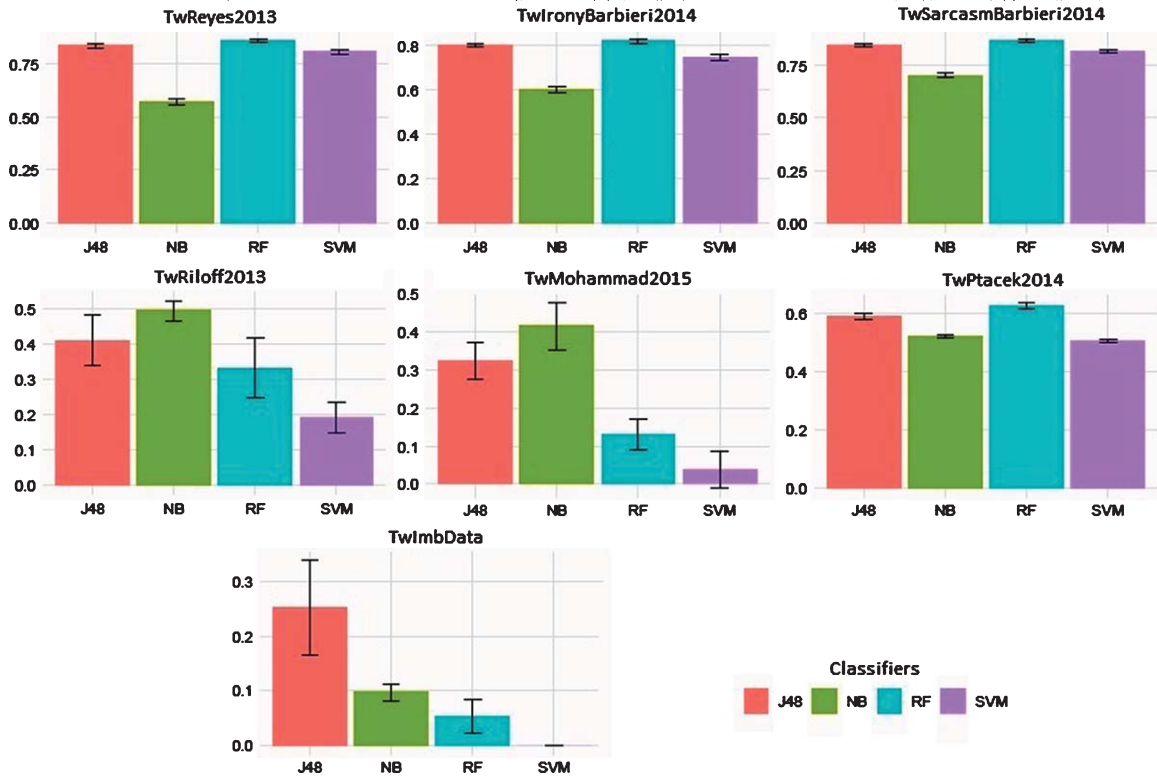


Fig. 8. Obtained results in F-score terms using the original distribution of the corpora.

### Balanced Accuracy

Figure 6 shows the results in terms of BAC over the ORIGINAL distribution of the corpora. Overall, the best results were obtained with RF, except in *TwRiloff2013*, where the highest rate of BAC was achieved by applying SVM.

The obtained results after applying imbalance treatment techniques in terms of BAC in most of the cases bring a drawback in the performance of the classifiers. However, in the case of *TwMohammad2015* and *TwImbData* (as it can be observed in Fig. 7) there is a positive impact when the three preprocessing techniques were applied. The performance of the classifiers after applying treatment methods on the *TwPtacek2014* shows the most significant drawback in terms of BAC compared with the ORIGINAL distribution. In terms of BAC, the most noticeable improvement was found over *TwImbData*.

### F-score

In Fig. 8 we present the obtained results in terms of F-score (it is the most widely applied evaluation metric in irony detection) when the experimental setting was applied using the ORIGINAL distribution.

As it can be noticed, the best performing algorithm in the ORIGINAL distribution was RF in most of the benchmark corpora, particularly in those that have been developed using the self-labeled approach. For what concerns the corpora involving a manual annotation process, the best performing classifier is NB. Concerning *TwImbData*, the J48 classifier obtains the highest results. As can be noticed, the F-score rates on *TwImbData* are lower than in the rest of the corpora reaching only 0.25 in F-score terms, while the highest score was near to 0.80 in *TwReyes2013* and *TwBarbieri2014*.

Figure 9 shows the obtained differences in terms of F-score. When applying the treatment techniques in *TwMohammad2015*, it is possible to improve the results of all classifiers, particularly of SVM. Regarding *TwRiloff2013*, the treatment techniques seem to have a positive impact on most of the experiments except when SMOTE was applied with NB and ROS with J48. Applying treatment techniques together with RF and SVM has a positive impact on the results involving *TwImbData*, while there is a drop in the results in both NB and J48. It is important to highlight that when RUS is used with J48 (the

553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575

576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599

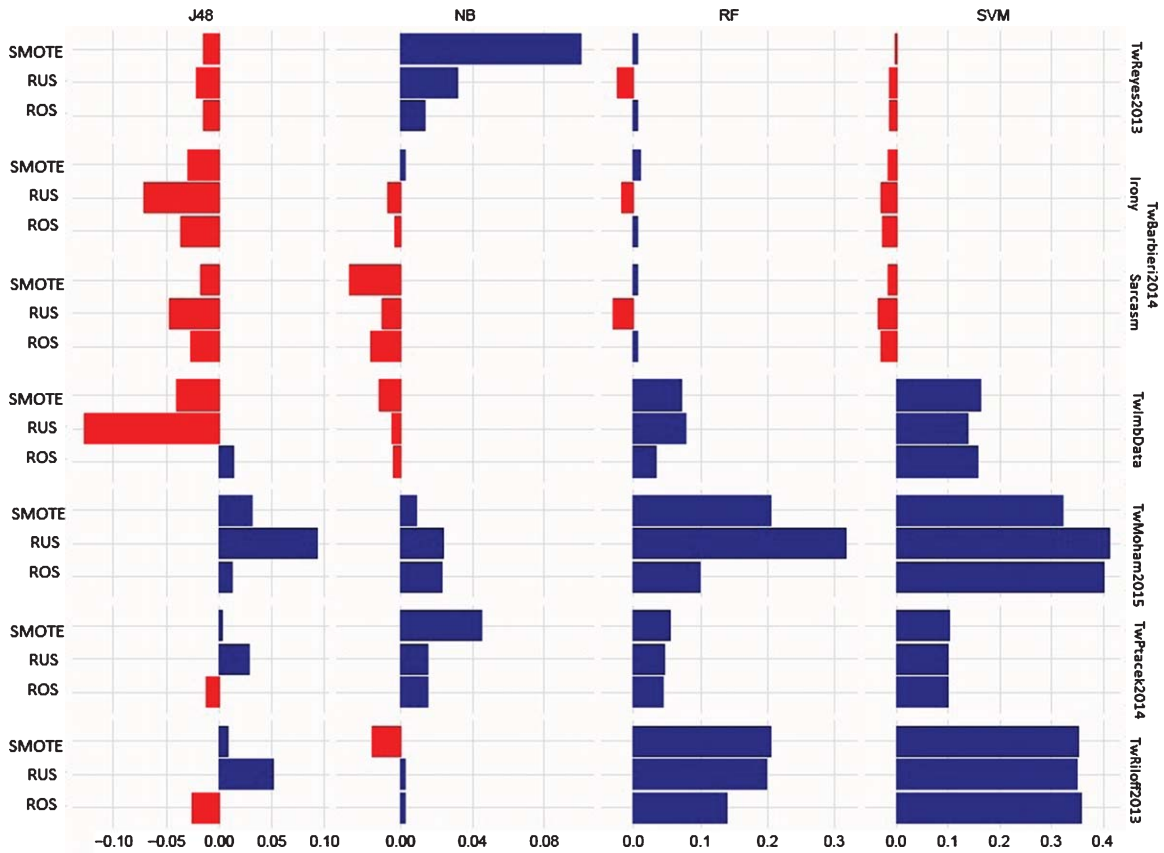


Fig. 9. Differences in terms of F-score with respect to the results of the ORIGINAL distribution after applying treatment techniques.

best performing classifier in the ORIGINAL distribution), its performance decreases. This could serve to validate the fact of the probability of losing useful information due to the nature of this treatment technique.

As already mentioned, F-score has been the most widely applied evaluation metric in the literature on irony detection. Therefore, by using this metric, it is possible to compare the performance of *emotIDM* when applying imbalance treatment techniques. Furthermore, unlike the rest of evaluation metrics used in this paper, it is possible to compare the obtained results against the related work.

Regarding the *TwReyes2013*, it is important to highlight that the experimental setting carried out in [17] for this dataset was different than in this paper. When *emotIDM* was evaluated over the aforementioned corpus, the authors considered a set of binary classifications between the ironic class and as negative instances each of the different subsets of tweets

(labeled with #education, #politics, etc.). For comparison purposes on the *TwReyes2013* the results reported in [13, 34] were considered; in these papers, the authors applied a similar setting than ours (i.e., the tweets belonging to the non-ironic classes were merged into a single class, and then a binary classification was carried out). The best performance on the ORIGINAL distribution outperforms the state of the art. Besides, when applying treatment techniques there are other classifiers obtaining better results than in the related work with a rate higher than 0.90 in F-score terms.

For what concerns to *TwBarbieri2014*, it is important to mention that there are not available results considering the same setting than in this paper, therefore it is not possible to compare the obtained results against the literature. In both subsets of *TwBarbieri2014*, the F-score rates are in some way similar to the ones obtained in *TwReyes2013*. Considering the ORIGINAL distribution, the best results were obtained with RF in both cases (irony and sarcasm).



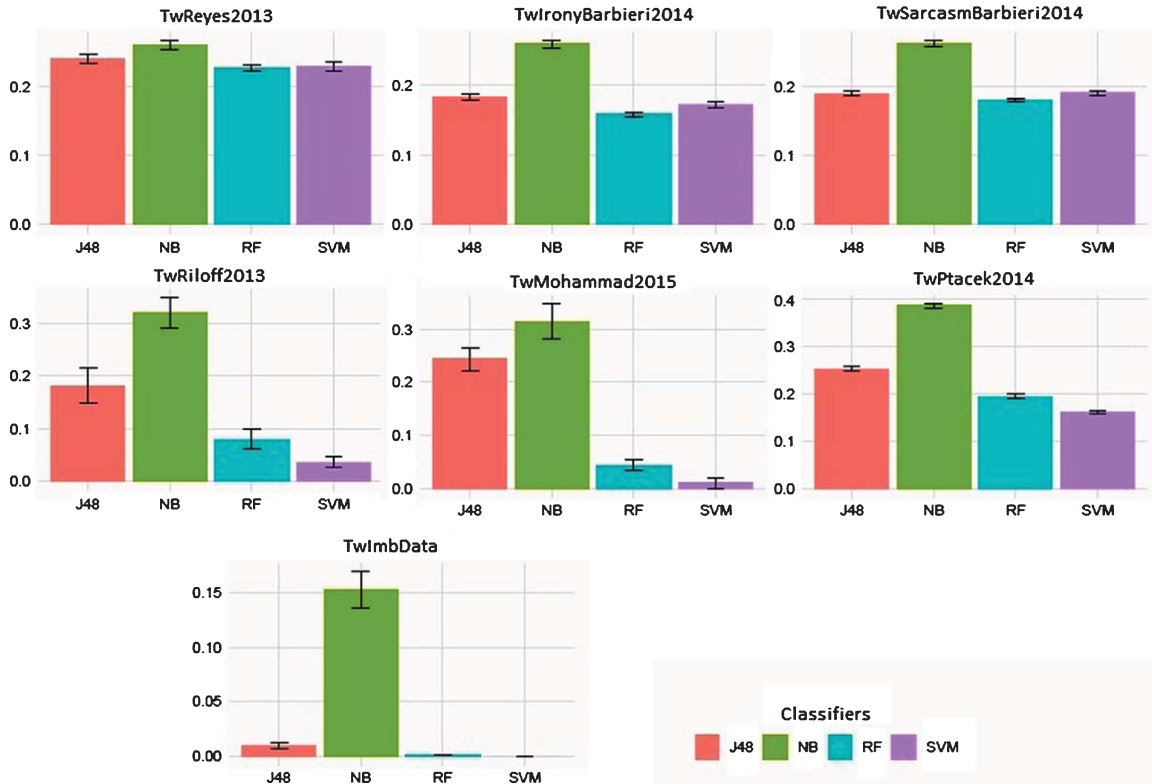


Fig. 10. Obtained results in Predictive Positives Percentage rate terms using the original distribution of the corpora.

Being the results in the *sarcasm-vs-non-sarcasm* experiments slightly better than in the case of *irony-vs-non-irony*.

Finally, in *TwImbData* the worst performing classifier in the ORIGINAL distribution is SVM. However, when the treatment techniques were applied, SVM emerges as the classifier having the best results.

### Predictive Positive Percentage

Finally, PPOS was used to show the percentage of instances classified as irony in each experiment. Figure 10 shows the performance in terms of Predictive Positive rate over the ORIGINAL distribution. The best performing classifier in terms of PPOS is NB. While in *TwMohammad2015*, *TwRiloff2013*, and *TwImbData*, SVM shows the worst results.

Figure 11 shows the obtained results after applying the imbalance treatment techniques. As can be observed, in all experiments there is an improvement in terms of PPOS. Overall, the highest results were achieved when applying RUS. While, the worst performance in terms of PPOS was obtained by using RF over *TwImbData* even after applying SMOTE or ROS.

## 4. Discussion of the results

In this section, we summarize the main findings of the experimental setting carried out applying different imbalance treatment techniques for addressing irony detection.

The RF classifier achieved the best results in terms of AUC, AUPR, and F-score in the case of self-labeled benchmark corpora. On the other hand, SVM showed the worst performance across the experiments, especially in those corpora with a high imbalanced class rate. In a similar fashion than in other domains, applying imbalance treatment techniques to the irony detection corpora before classifying with SVM, leads to an improvement in the performance, particularly in terms of Balanced Accuracy. However, there are some cases where applying imbalance treatment techniques provokes a drop in the performance of some classifiers. In terms of PPOS, it is possible to observe a positive impact on the performance of the classifiers, especially for *TwMohammad2015*, *TwRiloff2013*, *TwPtáček2014*, and *TwImbData*.

According to the results presented before, for each of the evaluation metrics, different imbalance treat-

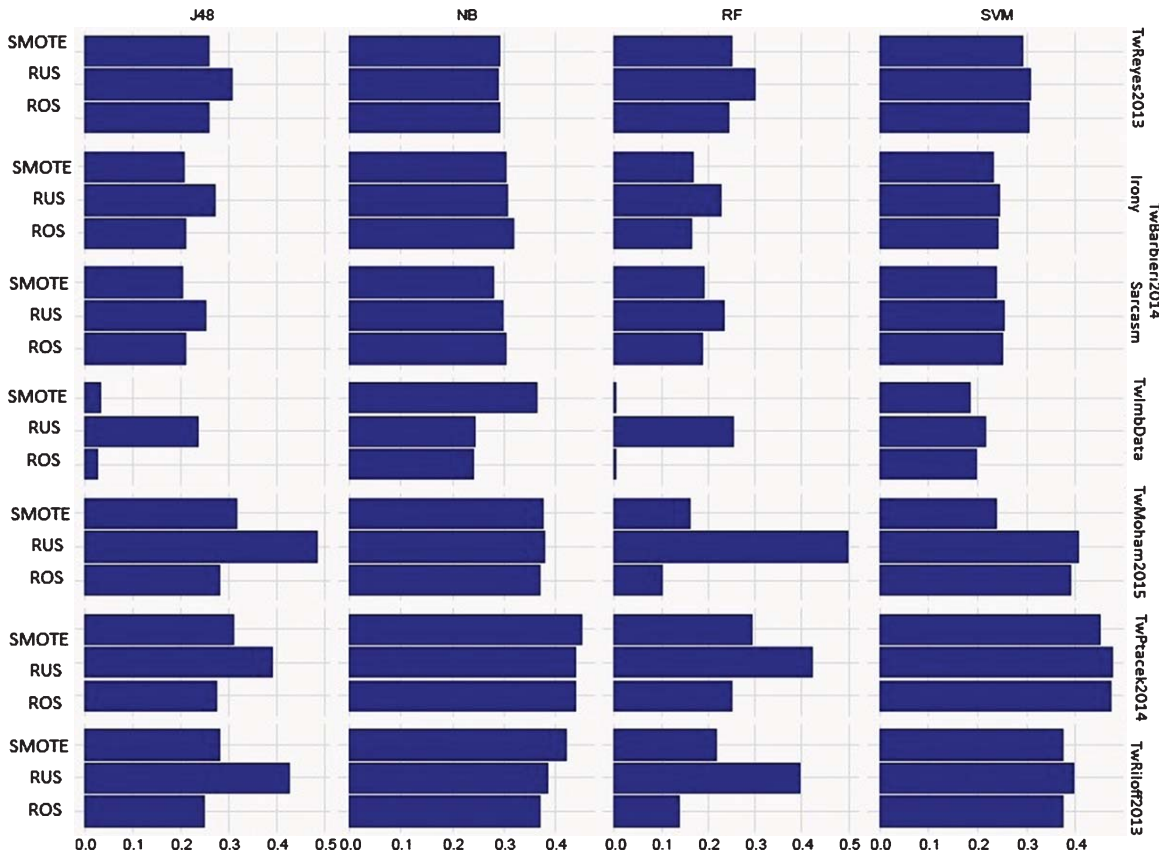


Fig. 11. Obtained results in terms of Predictive Positives Percentage rate considering the ORIGINAL distribution as well as applying treatment techniques.

689 ment techniques allow to improve the results of the  
 690 ORIGINAL distribution. SMOTE obtains the best  
 691 performance in terms of both AUPR and Balanced  
 692 Accuracy. Considering F-score, RUS is the method  
 693 allowing the best results. In terms of AUC, ROS  
 694 obtained the highest outcomes.

695 The corpora we used for experimental purposes  
 696 could be divided according to different aspects, for  
 697 example, considering the criteria used for retrieving  
 698 the data. The results in terms of PPOS in the ORIG-  
 699 INAL distribution seem to be higher when #sarcasm  
 700 is considered for retrieving data than in the case or  
 701 #irony.

702 Another aspect that can be considered within  
 703 the corpora we used concerns exploiting author's  
 704 self-labeled intention of being ironic (*TwReyes2013*,  
 705 *TwBarbieri2014*, *TwPtáček2014*, and *TwImbData*),  
 706 and the use of a manual annotation process  
 707 (*TwRiloff2013* and *TwMohammad2015*). In this case,  
 708 the results in terms of F-score in self-labeled cor-  
 709 pora are higher than in manually annotated data. This

710 could serve to validate the similar findings observed  
 711 in [16, 17] with reference to the impact of the corpora  
 712 construction methodology. However, on the other  
 713 hand, the most noticeable improvements on apply-  
 714 ing imbalance treatment techniques to compensate  
 715 the imbalance degree were achieved in those corpora  
 716 involving manual annotation.

717 Regarding the obtained results over *TwImbData*, it  
 718 is important to highlight that in all the evaluation met-  
 719 rics considered in this paper, there is a positive impact  
 720 on the performance of at least one of the classifiers.  
 721 *TwImbData* was developed having in mind to resem-  
 722 ble a realistic scenario where the difference between  
 723 ironic and non-ironic instances is very big. Therefore,  
 724 by improving the results over the ORIGINAL distri-  
 725 bution when the treatment methods were applied we  
 726 confirm the usefulness of using such techniques for  
 727 irony detection in imbalanced class scenarios.

728 Being irony a complex phenomenon, it is important  
 729 to assess the performance of different preprocessing  
 730 methods for compensating imbalance degree. As it

can be noticed, there is not a single imbalance treatment technique allowing to have the best performance across the evaluation metrics and corpora. This could be related to the nature of each method and also to the aim of the metrics.

As already mentioned, the most widely used evaluation metric in irony detection is F-score. Considering such a measure, the best performing technique was RUS, which serves to remove, in this case non-ironic samples. In our experimental setting, after applying imbalance treatment techniques both classes became balanced. Therefore, by applying RUS we are neither losing representative ironic instances nor generating synthetic instances.

Finally, it is important to highlight that the experiments were carried out only considering an irony detection model (*emoIDM*) relying mainly on affective features. It could be interesting to evaluate the performance of imbalance treatment techniques when ironic instances are represented by other kinds of features.

## 5. Conclusions

In this paper, we have evaluated the impact of class imbalance on detecting irony. We have performed several experiments over a set of Twitter corpora for irony detection covering different aspects such as corpora construction methodology and differences in data skew. Besides, we developed a set of irony corpora<sup>9</sup> aimed to resemble a more realistic scenario where the difference between the ironic and non-ironic class is very big. We employed *emoIDM*, an irony detection model based mainly on the presence of affective content. To the best of our knowledge, this is the first work in irony detection where a model for detecting such figurative language device is evaluated by considering many aspects related to the class imbalance problem.

In our research, we evaluated the performance of *emoIDM* together with a variety of classifiers when different imbalance treatment techniques were applied. Several metrics were used to compare the effectiveness of different classifiers and imbalance treatment techniques. Our results also allow us to compare the obtained results against those of the state of the art.

The main objective of this paper was to show that some treatment techniques can improve the per-

formance of classifiers dedicated to detect irony in Twitter particularly under an imbalanced class scenario. The results of this study indicate that the best performing imbalance treatment technique for addressing irony detection in imbalanced class scenarios depends on the evaluation metric used. However, considering the most widely used metric, i.e. F-score, the best performance was achieved by applying RUS.

We identified some directions for future work. It could be interesting to carry out some experiments using not only *data level approaches* (such as ROS, RUS, and SMOTE) but also *algorithm level approaches* (such as for example cost sensitive learning). Furthermore, experiments with other imbalance degree rates over the set of corpora used is part of the following steps of our research in irony detection in imbalanced class scenarios. On the other hand, it could be interesting to analyze the role of some of the data intrinsic characteristics described in [29] such as small disjuncts, lack of density and information as well as the overlapping between the classes on the irony detection corpora.

## Acknowledgments

The first author was funded by CONACYT project FC-2016/2410. Ronaldo Prati was supported by the São Paulo State (Brazil) research council FAPESP under project 2015/20606-6. Francisco Herrera was partially supported by the Spanish National Research Project TIN2017-89517-P. The work of Paolo Rosso was partially supported by the Spanish MICINN under the research project MISIMIS (PGC2018-096212-B-C31) and by the Generalitat Valenciana under the grant PROMETEO/2019/121.

## References

- [1] G. Abercrombie and D. Hovy, Putting Sarcasm Detection into Context: The Effects of Class Imbalance and Manual Labelling on Supervised Machine Classification of Twitter Conversations. In *Proceedings of the ACL 2016 Student Research Workshop*, pp. 107–113, Germany, August (2016).
- [2] L. Alba-Juez and S. Attardo, The Evaluative Palette of Verbal Irony, In Geoff Thompson and Laura Alba-Juez, editors, *Evaluation in Context*, pages 93–116. John Benjamins Publishing Company, Amsterdam/Philadelphia, (2014).
- [3] A. Ali, S.M. Shamsuddin and A.L. Ralescu, Classification with Class Imbalance Problem: A Review, *Intertional Jour-*

<sup>9</sup>The data will be released for research purposes.



- 826 *nal Of Advances In Soft Computing And Its Applications*  
827 **7**(3) (2015), 176–204.
- 828 [4] S. Attardo, Irony as Relevant Inappropriateness, In H. Col-  
829 ston and R. Gibbs, editors, *Irony in language and thought: A cognitive science reader*, pages 135–172. Lawrence Erl-  
830 baum, (2007).
- 831 [5] D. Bamman and N.A. Smith, Contextualized Sarcasm  
832 Detection on Twitter, In *Proceedings of the Ninth Inter-  
833 national Conference on Web and Social Media*, (2015), pp.  
834 574–577.
- 835 [6] F. Barbieri, H. Saggion and F. Ronzano, Modelling Sarcasm  
836 in Twitter, A Novel Approach. In *Proc. of the 5th Workshop  
837 on Computational Approaches to Subjectivity, Sentiment  
838 and Social Media Analysis*, pp. 50–58, USA, June (2014).  
839 ACL.
- 840 [7] G.E. de Almeida Prado Alves Batista, R.C. Prati and M.C.  
841 Monard, A Study of the Behavior of Several Methods for  
842 Balancing Machine Learning Training Data, *ACM Sigkdd  
843 Explorations Newsletter* **6**(1) (2004), 20–29.
- 844 [8] M.M. Bradley and P.J. Lang, Affective Norms for English  
845 Words (ANEW): Instruction Manual and Affective Ratings,  
846 Technical Report, Center for Research in Psychophysiology,  
847 University of Florida, Florida, (1999).
- 848 [9] E. Cambria, D. Olsher and D. Rajagopal, SenticNet  
849 3: A Common and Common-Sense Knowledge Base  
850 for Cognition-Driven Sentiment Analysis, In *Proceed-  
851 ings of AAAI Conference on Artificial Intelligence*, pages  
852 1515–1521, Canada, (2014). AAAI.
- 853 [10] A. Cervone, E.A. Stepanov, F. Celli and G. Riccardi, Irony  
854 Detection: From the Twittersphere to the News Space, In  
855 *Proceedings of Fourth Italian Conference on Computational  
856 Linguistics (CLiC-it 2017)*, (2017).
- 857 [11] N.V. Chawla, K.W. Bowyer, L.O. Hall and W. Philip  
858 Kegelmeyer, SMOTE: Synthetic Minority Oversampling  
859 Technique, *Journal of Artificial Intelligence Research* **16**  
860 (2002), 321–357.
- 861 [12] A. Fernández, S. García, M. Galar, R.C. Prati, B. Krawczyk  
862 and F. Herrera, Learning from imbalanced data sets,  
863 *Springer*, (2018).
- 864 [13] E. Fersini, F.A. Pozzi and E. Messina, Detecting Irony and  
865 Sarcasm in Microblogs: The Role of Expressive Signals and  
866 Ensemble Classifiers, In *International Conference on Data  
867 Science and Advanced Analytics*, (DSAA 2015), pages 1–8,  
868 France, (2015). IEEE Xplore Digital Library.
- 869 [14] H. Paul Grice, Logic and Conversation. In P. Cole and J.L.  
870 Morgan, editors, *Syntax and Semantics: Vol. 3 Speech Acts*,  
871 pages 41–58. Academic Press, (1975).
- 872 [15] H. He and E.A. Garcia. Learning from Imbalanced Data,  
873 *IEEE Transactions on Knowledge and Data Engineering*,  
874 **21**(9) (2009), 1263–1284.
- 875 [16] D.I.H. Fariás, M. M.-y. Gómez, H.J. Escalante, P. Rosso and  
876 V. Patti, A Knowledge-based Weighted KNN for Detecting  
877 Irony in Twitter, In *17th Mexican International Confer-  
878 ence on Artificial Intelligence (MICAI 2018)*, Mexico,  
879 (2018).
- 880 [17] D.I.H. Fariás, V. Patti and P. Rosso, Irony detection in Twit-  
881 ter: The Role of Affective Content, *ACM Trans Internet  
882 Technol* **16**(3) (2016), 19:1–19:24.
- 883 [18] D.I.H. Fariás and P. Rosso, Irony, Sarcasm, and Senti-  
884 ment Analysis. Chapter 7. In Federico A. Pozzi, Elisabetta  
885 Fersini, Enza Messina, and Bing Liu, editors, *Sentiment  
886 Analysis in Social Networks*, pages 113–127. Morgan Kauf-  
887 mann, (2016).
- 888 [19] M. Hu and B. Liu, Mining and Summarizing Customer  
889 Reviews. In *Proceedings of the Tenth ACM SIGKDD Inter-  
890 national Conference on Knowledge Discovery and Data  
891 Mining, KDD '04*, (2004), pages 168–177, USA, ACM.
- 892 [20] N. Japkowicz and S. Stephen, The Class Imbalance Prob-  
893 lem: A Systematic Study, *Intelligent Data Analysis* **6**(5)  
894 (2002), 429–449.
- 895 [21] A. Joshi, P. Bhattacharyya and M.J. Carman, Automatic  
896 Sarcasm Detection: A Survey. CoRR, abs/1602.03426,  
897 (2016).
- 898 [22] A. Joshi, V. Tripathi, K. Patel, P. Bhattacharyya and  
899 M.J. Carman, Are Word Embeddingbased Features Useful  
900 for Sarcasm Detection? In *Proc. of the 2016 Conference  
901 on Empirical Methods in Natural Language Processing*,  
902 Austin, Texas, USA, (2016), pages 1006–1011.
- 903 [23] J. Karoui, F. Benamara, V. Moriceau, N. Aussenac-Gilles  
904 and L. Hadrich-Belguith, Towards a Contextual Pragmatic  
905 Model to Detect Irony in Tweets. In *Proc. of the 53rd Annual  
906 Meeting of the Association for Computational Linguistics  
907 and the 7th International Joint Conference on Natural Lan-  
908 guage Processing (Volume 2: Short Papers)*, pages 644–650,  
909 China, (2015). ACL.
- 910 [24] A. Khattri, A. Joshi, P. Bhattacharyya and M. Carman,  
911 Your Sentiment Precedes You: Using an author’s Historical  
912 Tweets to Predict Sarcasm. In *Proceedings of the 6th  
913 Workshop on Computational Approaches to Subjectivity,  
914 Sentiment and Social Media Analysis*, pages 25–30, Por-  
915 tugal, September (2015). ACL.
- 916 [25] S. Kumon-Nakamura and S. Glucksberg, How About  
917 Another Piece of Pie: The Allusional Pretense Theory of  
918 Discourse Irony, *Journal of Experimental Psychology: Gen-  
919 eral* **124**(1) (1995), 3.
- 920 [26] F. Kunneman, C. Liebrecht, M. van Mulken and A. van  
921 den Bosch, Signaling Sarcasm: From Hyperbole to Hash-  
922 tag, *Information Processing & Management* **51**(4) (2015),  
923 500–509.
- 924 [27] C. Liebrecht, F. Kunneman and A. Van den Bosch, The  
925 Perfect Solution for Detecting Sarcasm in Tweets #not.  
926 In *Proceedings of the 4th Workshop on Computational  
927 Approaches to Subjectivity, Sentiment and Social Media  
928 Analysis*, pages 29–37, USA, June (2013). ACL.
- 929 [28] P. Liu, W. Chen, G. Ou, T. Wang, D. Yang and K.  
930 Lei, Sarcasm Detection in Social Media Based on Imbal-  
931 anced Classification. In Feifei Li, Guoliang Li, Seungwon  
932 Hwang, Bin Yao, and Zhenjie Zhang, editors, *Proceedings  
933 of the Web-Age Information Management: 15th Internation-  
934 al Conference*, pages 459–471, China, (2014). Springer  
935 International Publishing.
- 936 [29] V. López, A. Fernández, S. García, V. Palade and F. Herrera,  
937 An Insight into Classification with Imbalanced Data: Empir-  
938 ical Results and Current Trends on Using Data Intrinsic  
939 Characteristics, *Information Sciences* **250** (2013), 113–141.
- 940 [30] S. McDonald, Neuropsychological Studies of Sarcasm. In  
941 H. Colston and R. Gibbs, editors, *Irony in language and  
942 Thought: A Cognitive Science Reader*, pages 217–230.  
943 Lawrence Erlbaum, (2007).
- 944 [31] S.M. Mohammad and P.D. Turney, Crowdsourcing a  
945 Word–emotion Association Lexicon, *Computational Intel-  
946 ligence* **29**(3) (2013), 436–465.
- 947 [32] S.M. Mohammad, X. Zhu, S. Kiritchenko and J. Mar-  
948 tin, Sentiment, Emotion, Purpose, and Style in Electoral  
949 Tweets, *Information Processing & Management* **51**(4)  
950 (2015), 480–499.
- 951 [33] F.Å. Nielsen, A New ANEW: Evaluation of a Word List  
952 for Sentiment Analysis in Microblogs, In *Proceedings of  
953 the ESWC2011 Workshop on 'Making Sense of Microposts':  
954 Big things come in small packages*, volume **718**  
955

- 956 of CEUR Workshop Proceedings, pages 93–98, Greece,  
957 (2011). CEUR-WS.org.
- [34] D. Nozza, E. Fersini and E. Messina, Unsupervised Irony  
958 Detection: A Probabilistic Model with Word Embeddings,  
959 In *Proceedings of the 8th International Joint Conference on*  
960 *Knowledge Discovery, Knowledge Engineering and Knowl-*  
961 *edge Management*, pages (2016), 68–76.
- [35] J.W. Pennebaker, M.E. Francis and R.J. Booth, Linguistic  
962 Inquiry and Word Count: LIWC 2001, *Mahway: Lawrence*  
963 *Erlbaum Associates* **71** (2001), 2–23.
- [36] S. Poria, E. Cambria, D. Hazarika and P. Viji, A Deeper Look  
964 into Sarcastic Tweets Using Deep Convolutional Neural  
965 Networks, In *Proceedings of COLING 2016, the 26th Inter-*  
966 *national Conference on Computational Linguistics*, page  
967 1601–1612, Japan, December (2016).
- [37] S. Poria, A. Gelbukh, A. Hussain, N. Howard, D. Das and S.  
968 Bandyopadhyay, Enhanced SenticNet with Affective Labels  
969 for Concept-Based Opinion Mining, *IEEE Intelligent Sys-*  
970 *tems* **28**(2) (2013), 31–38.
- [38] R.C. Prati, Gustavo Enrique de Almeida Prado Alves  
971 Batista, and Diego F. Silva. Class Imbalance Revisited:  
972 A New Experimental setup to Assess the Performance of  
973 Treatment Methods, *Knowledge and Information Systems*  
974 **45**(1) (2015), 247–270.
- [39] T. Ptáček, I. Habernal and J. Hong, Sarcasm Detection on  
975 Czech and English Twitter, In *Proceedings of COLING*  
976 *2014, the 25th International Conference on Computational*  
977 *Linguistics*, pages 213–223, Ireland, August (2014). Dublin  
978 City University and ACL.
- [40] A. Reyes, P. Rosso and T. Veale, A Multidimensional  
979 Approach for Detecting Irony in Twitter, *Language*  
980 *Resources and Evaluation* **47**(1) (2013), 239–268.
- [41] E. Riloff, A. Qadir, P. Surve, L. De Silva, N. Gilbert and  
981 R. Huang, Sarcasm as Contrast between a Positive Sentiment  
982 and Negative Situation, In *Proceedings of the 2013*  
983 *Conference on Empirical Methods in Natural Language*  
984 *Processing*, (EMNLP 2013), pages 704–714. USA, October  
985 (2013). ACL.
- [42] E. Sulis, D.I.H. Fariás, P. Rosso, V. Patti and G. Ruffo, Fig-  
986 urative Messages and Affect in Twitter: Differences between  
987 #irony, #sarcasm and #not, *Knowledge-Based Systems* **108**  
988 (2016), 132–143.
- [43] A. Utsumi, Verbal Irony as Implicit Display of Ironic Envi-  
989 ronment: Distinguishing Ironic Utterances from Nonirony,  
990 *Journal of Pragmatics* **32**(12) (2000), 1777–1806.
- [44] B.C. Wallace, D.K. Choe and E. Charniak, Sparse, Con-  
991 textually Informed Models for Irony Detection: Exploiting  
992 User Communities, *Entities and Sentiment*. In *Proc. of the*  
993 *53rd Annual Meeting of the Association for Computational*  
994 *Linguistics and the 7th International Joint Conference on*  
995 *Natural Language Processing (Volume 1: Long Papers)*,  
996 pages 1035–1044, China, July 2015. ACL.
- [45] A.P. Wang. #Irony or #Sarcasm — A Quantitative and  
997 Qualitative Study Based on Twitter. In *Proceedings of the*  
998 *PACLIC: the 27th Pacific Asia Conference on Language,*  
999 *Information, and Computation*, pages 349–356, Taiwan,  
1000 (2013).
- [46] C. Whissell, Using the Revised Dictionary of Affect in Lan-  
1001 guage to Quantify the Emotional Undertones of Samples of  
1002 Natural Languages, *Psychological Reports* **2**(105) (2009),  
1003 509–521.
- [47] D. Wilson and D. Sperber, On Verbal Irony, *Lingua* **87**(1-2)  
1004 (1992), 53–76.