



## **ADVCOMP 2020**

The Fourteenth International Conference on Advanced Engineering Computing  
and Applications in Sciences

ISBN: 978-1-61208-812-9

October 25 - 29, 2020

### **ADVCOMP 2020 Editors**

Claus-Peter Rückemann, Westfälische Wilhelms-Universität Münster (WWU) /

DIMF / Leibniz Universität Hannover, Germany

Tim vor der Brück, FFHS, Lucerne University of Applied Sciences and Arts,

Switzerland

# ADVCOMP 2020

## Forward

The Fourteenth International Conference on Advanced Engineering Computing and Applications in Sciences (ADVCOMP 2020), held on October 22-29, 2020, continued a series of events meant to bring together researchers from the academia and practitioners from the industry in order to address fundamentals of advanced scientific computing and specific mechanisms and algorithms for particular sciences.

With the advent of high performance computing environments, virtualization, distributed and parallel computing, as well as the increasing memory, storage and computational power, processing particularly complex scientific applications and voluminous data is more affordable. With the current computing software, hardware and distributed platforms effective use of advanced computing techniques is more achievable.

The conference provided a forum where researchers were able to present recent research results and new research problems and directions related to them. The conference sought contributions presenting novel research in all aspects of new scientific methods for computing and hybrid methods for computing optimization, as well as advanced algorithms and computational procedures, software and hardware solutions dealing with specific domains of science.

The conference had the following tracks:

- Computing applications in science
- Computing mechanisms and methods
- Multidisciplinary Mobile and Web Applications in Modern Life

We take here the opportunity to warmly thank all the members of the ADVCOMP 2020 technical program committee, as well as all the reviewers. The creation of such a high quality conference program would not have been possible without their involvement. We also kindly thank all the authors that dedicated much of their time and effort to contribute to ADVCOMP 2020. We truly believe that, thanks to all these efforts, the final conference program consisted of top quality contributions.

We also gratefully thank the members of the ADVCOMP 2020 organizing committee for their help in handling the logistics and for their work that made this professional meeting a success.

We hope that ADVCOMP 2020 was a successful international forum for the exchange of ideas and results between academia and industry and to promote further progress in the field of engineering computing and applications in sciences.

### **ADVCOMP 2020 General Chair**

Sandra Sendra, Universitat Politecnica de Valencia, Universidad de Granada, Spain

### **ADVCOMP 2020 Steering Committee**

Juha Röning, University of Oulu, Finland

Hans-Joachim Bungartz, TUM, Germany

Andreas Rausch, Technische Universität Clausthal, Germany

Dean Vucinic, Vrije Universiteit Brussel (VUB), Belgium | FERIT, Croatia,

Marcin Hojny, AGH University of Science and Technology, Poland

Alice E. Koniges, University of Hawai'i at Mānoa, USA

Alfred Geiger, T-Systems Information Services GmbH, Germany

### **ADVCOMP 2020 Publicity Chair**

Daniel Andoni Basterrechea, Universitat Politecnica de Valencia, Spain

## **ADVCOMP 2020**

### **COMMITTEE**

#### **ADVCOMP General Chair**

Sandra Sendra, Universitat Politecnica de Valencia, Universidad de Granada, Spain

#### **ADVCOMP Steering Committee**

Dean Vucinic, Vrije Universiteit Brussel (VUB), Belgium | FERIT, Croatia,

Alfred Geiger, T-Systems Information Services GmbH, Germany

Hans-Joachim Bungartz, TUM, Germany

Andreas , Technische Universität Clausthal, Germany

Juha Röning, University of Oulu, Finland

Marcin Hojny, AGH University of Science and Technology, Poland

Alice E. Koniges, University of Hawai'i at Mānoa, USA

#### **ADVCOMP 2020 Publicity Chair**

Daniel Andoni Basterrechea, Universitat Politecnica de Valencia, Spain

#### **ADVCOMP 2020 Technical Program Committee**

Mohd Helmy Abd Wahab, Universiti Tun Hussein Onn Malaysia, Malaysia

Waleed H. Abdulla, University of Auckland, New Zealand

José Abellán, Catholic University of Murcia, Spain

Mohamed Riduan Abid, Alakhawayn University, Morocco

Francisco Airton Silva, Federal University of Piauí, Brazil

M. Azeem Akbar, Nanjing University of Aeronautics and Astronautics, China

Haifa Alharthi, Saudi Electronic University, Saudi Arabia

Sónia Maria Almeida da Luz, Polytechnic Institute of Leiria - School of Technology and Management, Portugal

Madyan Alsenwi, Kyung Hee University, Global Campus, South Korea

Mohamed E. Aly, California State Polytechnic University, Pomona, USA

Daniel Andresen, Kansas State University, USA

Alberto Antonietti, Politecnico di Milano / University of Pavia, Italy

Ehsan Atoofian, Lakehead University, Canada

Vadim Azhmyakov, Universidad EAFIT, Medellin, Republic of Colombia

Carlos Becker Westphall, University of Santa Catarina, Brazil

Raoudha Ben Djemaa, ISITCOM | University of Sousse, Tunisia

Rudolf Berrendorf, Bonn-Rhein-Sieg University, Germany

Estêvão Bissoli Saleme, Federal University of Espírito Santo, Brazil

Alessandro Borri, CNR-IASI Biomathematics Laboratory, Rome, Italy

David Bouck-Standen, Kingsbridge Research Center, UK

Hans-Joachim Bungartz, TUM, Germany

Xiao-Chuan Cai, University of Colorado Boulder, USA

Juan Vicente Capella Hernández, Universitat Politècnica de València, Spain  
Graziana Cavone, Polytechnic of Bari, Italy  
Mete Celik, Erciyes University, Turkey  
Jinyuan Chen, Louisiana Tech University, USA  
Rangeen Basu Roy Chowdhury, Intel Corporation, USA  
Robert Clay, Sandia National Labs, USA  
Vassilios V. Dimakopoulos, University of Ioannina, Greece  
Inês Domingues, IPO Porto Research Centre (CI-IPOP), Portugal  
Shi Dong, Northeastern University, Boston, USA  
Maha Elarbi, University of Tunis, Tunisia  
Javier Fabra, Universidad de Zaragoza, Spain  
Akemi Galvez, University of Cantabria, Spain / Toho University, Japan  
Félix J. García Clemente, University of Murcia, Spain  
Leonardo Garrido, Tecnológico de Monterrey, Mexico  
Alfred Geiger, T-Systems Information Services GmbH, Germany  
Tong Geng, Boston University, USA  
Jing Gong, KTH Royal Institute of Technology, Sweden  
Teofilo Gonzalez, UC Santa Barbara, USA  
Bernard Grabot, LGP-ENIT, France  
Maki Habib, American University in Cairo, Egypt  
Marcin Hojny, AGH University of Science and Technology, Poland  
Wladyslaw Homenda, Warsaw University of Technology, Poland  
Tzung-Pei Hong, National University of Kaohsiung, Taiwan  
Mehdi Hosseinzadeh, Washington University in St. Louis, USA  
Paul Humphreys, Ulster University | Ulster University Business School, UK  
Andres Iglesias, University of Cantabria, Spain / Toho University, Japan  
Hiroshi Ishikawa, Tokyo Metropolitan University, Japan  
Attila Kertesz, University of Szeged, Hungary  
Alice E. Koniges, University of Hawai'i at Mānoa, USA  
V M Krushnarao Kottedda, University of Wyoming, USA  
Seyong Lee, Oak Ridge National Laboratory, USA  
Maurizio Leotta, University of Genova, Italy  
Clement Leung, Chinese University of Hong Kong, Shenzhen, China  
Yiu-Wing Leung, Hong Kong Baptist University, Hong Kong  
Yiheng Liang, Bridgewater State University, USA  
Stephane Maag, Telecom SudParis, France  
Elbert E. N. Macau, Federal University of Sao Paulo - UNIFESP at Sao Jose dos Campos, Brazil  
Marcin Markowski, Wroclaw University of Science and Technology, Poland  
Mirko Marras, University of Cagliari, Italy  
Simon Mille, Universitat Pompeu Fabra, Spain  
Mohamed Wiem Mkaouer, Rochester Institute of Technology, USA  
Sébastien Monnet, Savoie Mont Blanc University (USMB), France  
Shana Moothedath, University of Washington, Seattle, USA  
Laurent Nana, University of Brest, France  
Ehsan Nekouei, City University of Hong Kong, Hong Kong  
Joanna Isabelle Olszewska, University of West Scotland, UK  
Diego P. Ruiz, University of Granada, Spain  
Marcin Paprzycki, Systems Research Institute | Polish Academy of Sciences, Poland

Prantosh Kumar Paul, Raiganj University, India  
Damien Pellier, Université Grenoble Alpes, France  
Sonia Pérez-Díaz, University of Alcalá, Spain  
Antonio Petitti, Institute of Intelligent Industrial Systems and Technologies for Advanced Manufacturing (STIIMA) - National Research Council of Italy (CNR), Italy  
Tamas Pflanzner, University of Szeged, Hungary  
Agostino Poggi, Università degli Studi di Parma, Italy  
Andreas Rausch, Technische Universität Clausthal, Germany  
Carlos Reaño, Queen's University Belfast, UK  
Michele Risi, University of Salerno, Italy  
Michele Roccotelli, Politecnico di Bari, Italy  
Ivan Rodero, Rutgers University, USA  
Juha Röning, University of Oulu, Finland  
Julio Sahuquillo, Universitat Politècnica de València, Spain  
Subhash Saini, NASA, USA  
Hamed Sarvari, George Mason University, USA  
Alireza Shahrabi, Glasgow Caledonian University, Scotland, UK  
Justin Shi, Temple University, USA  
Philip Shilane, Dell EMC, USA  
Costas Vassilakis, University of the Peloponnese, Greece  
Flavien Vernier, LISTIC – Savoie University, France  
Dean Vucinic, Vrije Universiteit Brussel (VUB), Belgium / FERIT, Croatia  
Adriano V. Werhli, Universidade Federal do Rio Grande - FURG, Brazil  
Gabriel Wittum, Goethe University Frankfurt, Germany  
Mudasser F. Wyne, National University, USA  
Cong-Cong Xing, Nicholls State University, USA  
Feng Yan, University of Nevada, Reno, USA  
Carolina Yukari Veludo Watanabe, Federal University of Rondônia, Brazil  
Michael Zapf, Technische Hochschule Nürnberg Georg Simon Ohm (University of Applied Sciences Nuremberg), Germany  
Vesna Zeljkovic, Lincoln University, USA  
Ruochen Zeng, NXP Semiconductors, USA  
Penghui Zhang, Arizona State University, USA  
Qian Zhang, Liverpool John Moores University, UK

## Copyright Information

For your reference, this is the text governing the copyright release for material published by IARIA.

The copyright release is a transfer of publication rights, which allows IARIA and its partners to drive the dissemination of the published material. This allows IARIA to give articles increased visibility via distribution, inclusion in libraries, and arrangements for submission to indexes.

I, the undersigned, declare that the article is original, and that I represent the authors of this article in the copyright release matters. If this work has been done as work-for-hire, I have obtained all necessary clearances to execute a copyright release. I hereby irrevocably transfer exclusive copyright for this material to IARIA. I give IARIA permission to reproduce the work in any media format such as, but not limited to, print, digital, or electronic. I give IARIA permission to distribute the materials without restriction to any institutions or individuals. I give IARIA permission to submit the work for inclusion in article repositories as IARIA sees fit.

I, the undersigned, declare that to the best of my knowledge, the article does not contain libelous or otherwise unlawful contents or invading the right of privacy or infringing on a proprietary right.

Following the copyright release, any circulated version of the article must bear the copyright notice and any header and footer information that IARIA applies to the published article.

IARIA grants royalty-free permission to the authors to disseminate the work, under the above provisions, for any academic, commercial, or industrial use. IARIA grants royalty-free permission to any individuals or institutions to make the article available electronically, online, or in print.

IARIA acknowledges that rights to any algorithm, process, procedure, apparatus, or articles of manufacture remain with the authors and their employers.

I, the undersigned, understand that IARIA will not be liable, in contract, tort (including, without limitation, negligence), pre-contract or other representations (other than fraudulent misrepresentations) or otherwise in connection with the publication of my work.

Exception to the above is made for work-for-hire performed while employed by the government. In that case, copyright to the material remains with the said government. The rightful owners (authors and government entity) grant unlimited and unrestricted permission to IARIA, IARIA's contractors, and IARIA's partners to further distribute the work.

## Table of Contents

The F-Measure Paradox <i>Tim vor der Bruck</i>	1
Detecting Suicide Risk Through Twitter <i>Javier Fabra, Ana Belen Martinez-Martinez, Yolanda Lopez-Del-Hoyo, Maria Cruz Perez-Yus, and Barbara Olivan-Blazquez</i>	6
Adaptive Multimedia Indexing Using Naive Bayes Classification <i>R. F. Ma, Clement Leung, and Jiayan Zhang</i>	12



# The F-Measure Paradox

Tim vor der Brück

Fernfachhochschule Schweiz (FFHS)  
 Brig, Switzerland  
 Email: tim.vorderbrueck@ffhs.ch

**Abstract**—Paradoxes have raised a lot of interest in mathematics and computer science. What fascinates people about them is that such a paradox contains a self-contradictory statement that dissents with usual beliefs and expectations. The range of discovered paradoxes is long. One of the most famous is probably the proposition of Russell that states that no set can exist that contains all sets that do not contain itself as a subset. The paradox arises in the proof, where it is shown that such a set must contain itself if and only if it does not contain itself. In this paper, we derive a paradox about the F-measure, one of the most important metrics in machine learning. The contribution of this paper is two-fold. On the one hand, we investigate typical properties of the F-Measure, on the other hand, we show that they are contradictory and therefore constitute a paradox, to several properties of the harmonic mean, where the F-Measure is a special case of.

**Keywords**—F-Measure; paradox; precision; recall; NaN.

## I. INTRODUCTION

The word *paradox* originates from Greece and is composed of the word *para* (beyond) and *doxa* (opinion). A paradox contains a self-contradictory statement and dissents with people’s beliefs and expectations [1]. It often does not have a direct practical use case but it gives theoretical insights and helps to understand certain problems better. Especially in the area of mathematics, there is a large amount of identified paradoxes. A quite well-known paradox is the proposition of Russel.

*Russel’s Paradox*: This proposition [2] claims that no set can exist that contains all sets that do not contain themselves and nothing more. The proof is done by contradiction. Let us assume such a set would exist. Then exactly one of the following propositions must be true about this set:

- This set contains itself. This is not possible since this set only contains sets that do not contain themselves.
- This set does not contain itself. Then per definition, this set must contain itself, which is a contradiction.

Since both cases lead to a contradiction, such a set cannot exist.

*Banach-Tarski Paradox*: Another well-known paradox from mathematics is the so-called Banach-Tarski-Paradox [3] that claims that a sphere can be decomposed and put together afterward in such a way that one has obtained two spheres of the same volume as the original sphere. Thus, one of the spheres was seemingly created out of nothing. This paradox is based on the principle that some concepts of mathematics cannot be transferred into reality.

		Prediction outcome		total
		p	n	
Actual value	p'	8	2	P'
	n'	12	9978	N'
total		P	N	

Figure 1. Example of confusion matrix for an imbalanced class distribution.

*Stein’s Paradox*: Normally, the expected value is best approximated by the average value, since the average value is actually its best unbiased estimator. Stein’s paradox [4] states that, if several expected values of the same type are to be determined (like batting statistics for a collection of baseball players), the isolated averages are no longer the best choice. Instead, all the estimates should be determined jointly by shifting the individual estimates in direction of the overall cross-estimate average.

*Accuracy Paradox*: Related to Data Science is the so-called accuracy paradox [5][6]. It states that when comparing two classification methods, the one with the lower accuracy can have in fact higher predictory capability. This phenomenon usually occurs in the case of highly imbalanced class distributions. Consider for example a very infrequent event like a rare disease that only shows up for around 0.1% of the cases. Let us assume, we have a method that can detect 40% of the events correctly and its precision is 80%. So, its confusion matrix could look like the one in Figure 1, where the columns denote the predicted and the rows the actual values. The obtained accuracy of this method would then amount to  $9986/10000=0.9986$  while predicting always the majority class (event not occurring), which has in fact no predictive power, would achieve an accuracy of 0.999.

In this paper, we will first derive several general statements about the harmonic mean of two variables. Afterward, we will proof that these statements are indeed incorrect for the  $F_1$ -score, which is a special case of the harmonic mean, in partic-

ular, it is the harmonic mean of precision and recall. Finally, we will analyze the reasons for this paradox and investigate the consequences for mathematical proofs in general.

The remainder of this paper is organized as follows. In Section II, we define the general harmonic mean and show several of its principal properties. Section III gives an overview of the  $F_1$ -score, which is the harmonic mean of precision and recall. In the next section (Section IV), the paradox of the  $F_1$ -score is described. The findings and the cause of this paradox are discussed in the next Section V. Finally, the paper concludes with Section VI, which summarizes the obtained results.

## II. HARMONIC MEAN

The harmonic mean  $H(a, b)$  of two values  $a$  and  $b$  is the Hoelder-mean with coefficient -1 and is formally given by [7]:

$$H(a, b) = \left( \frac{a^{-1} + b^{-1}}{2} \right)^{-1} = \frac{2}{\frac{1}{a} + \frac{1}{b}} \quad (1)$$

with  $a, b \in \mathbb{R}$  (or  $\mathbb{C}$ ). An alternative and simpler formulation is:

$$H(a, b) = \frac{2ab}{a+b} \quad (2)$$

In contrast to the arithmetic mean, the harmonic mean is a rather pessimistic mean that is drawn in direction to the minimum of both arguments. Note that it can only be applied to argument values of identical signs [7]. To see this, consider the following example:

$$H(-2, 4) = \frac{2(-2) \cdot 4}{-2+4} = \frac{-16}{2} = -8 \notin [-2, 4]$$

Since -8 is not located between -2 and 4, it cannot possibly constitute any mean of those values.

Consider now the following two propositions (1+2):

$$\begin{aligned} 1. a = 0 &\Rightarrow H(a, b) = 0 \\ 2. H(a, b) = 0 &\Rightarrow a = 0 \end{aligned} \quad (3)$$

Note that without limitation of generality  $a=0$  can be replaced by  $b=0$  due to the symmetry of the harmonic mean.

It is not difficult to show that the first statement is true and the second false.

*Proof:* Let us first have a look at proposition 1. The following two cases can be discerned:  $b \neq 0$  and  $b = 0$ . First, we consider the case  $b \neq 0$ . Plugging in  $a = 0$  in formula 2 results in:

$$H(a, b) = \frac{2 \cdot 0 \cdot b}{0 + b} = \frac{0}{b} = 0 \quad (4)$$

Now consider  $a = 0, b = 0$ . Plugging both values into  $H(a, b)$  results in an expression  $\frac{0}{0}$ , which is not defined. Let us, however, look at the behavior of  $H(a, b)$  for  $a$  and  $b$  approaching zero using formula 1. Since the sign of  $a$  and  $b$  must coincide, we get:

$$\lim_{a, b \rightarrow 0} H(a, b) = \frac{2}{\frac{1}{a} + \frac{1}{b}} = \frac{2}{\infty} = 0 \quad (5)$$

Therefore, it is a reasonable approach to define  $H(0, 0)$  as 0, to which we henceforth abide.

Proposition 2 is straight-forward to show by the following counterexample:  $H(2, 0) = 0$  but  $2 \neq 0$ . ■

One can also draw some conclusions, under which conditions the harmonic mean  $H$  and one of its input arguments have to coincide. In particular, assuming  $a$  is not diminishing, then from the fact that  $a$  and  $H$  coincide one can infer that  $b$  must also assume their common value. The opposite, however, is false.

Formally, the first proposition (proposition 3) is true and the second (proposition 4) is false:

$$\begin{aligned} 3. a \neq 0 \wedge a = H(a, b) &\Rightarrow b = H(a, b) \\ 4. a \neq 0 \wedge b = H(a, b) &\Rightarrow a = H(a, b) \end{aligned} \quad (6)$$

*Proof of Proposition 3:* From the definition of the harmonic mean, it follows that:  $H(a, b) = \frac{2ab}{a+b}$ . Since  $H(a, b)$  equals  $a$ , we can plugin  $a$  on the left-hand side:  $a = \frac{2ab}{a+b}$ . Since  $a \neq 0$ , both sides can be divided by  $a$ :  $1 = \frac{2b}{a+b}$ . Afterward, we multiply both sides by  $a + b$ :  $a + b = 2b$ . By subtracting  $b$  from both sides one finally obtains:  $a = b$ . ■

The opposite direction (proposition 4) can be shown by contraction, let  $a=1, b=0=H(a,b)$ , then herewith it follows that  $a \neq H(a, b)$ .

## III. $F_1$ -SCORE

The  $F_1$ -Score is the harmonic mean of precision and recall, where precision is the percentage of predicted positive events that are indeed positive, while recall is the percentage of positive events that are actually correctly detected by the algorithm [8]. All three measures originated from the area of information retrieval but quickly spread into other areas of machine learning too. Let TP be the true positives, i.e., the number of positive events that were correctly classified by the algorithm, FP the number of negative events that were actually classified as positive, and FN the number of positive events that were misclassified as negative. Then precision (prec), recall (rec), and F-measure are formally defined as follows:

$$\begin{aligned} prec &= \frac{TP}{TP + FP} \\ rec &= \frac{TP}{TP + FN} \\ F_1(prec, rec) &= H(prec, rec) \\ &= \frac{2prec \cdot rec}{prec + rec} \end{aligned} \quad (7)$$

Note that recall or precision can potentially be undefined. Consider, for example, that the positive class never shows up in the evaluation data. In this case, TP and FN assume both zero, which results in an undefined recall value. Similarly, if the positive class is never predicted, the precision is left undefined. Analogously to the definition of floating point numbers, we use the expression  $NaN$  to denote an undefined value, which

stands for *Not a Number*. We also define arithmetic on  $NaN$  in the following way by following the Bochvar extension [9]. Let  $a \in \mathbb{R} \cup \{NaN\}$  be arbitrarily chosen, then:

$$\begin{aligned} a \cdot NaN &= NaN \\ a + NaN &= NaN \\ a - NaN &= NaN \\ \frac{a}{NaN} &= NaN \end{aligned} \tag{8}$$

As one can easily perceive, if at least one of the operator arguments assumes  $NaN$ , then also the result is  $NaN$ . Therefore,  $NaN$  is also called an absorbing element. Regarding the algebraic structure,  $\mathbb{R} \cup \{NaN\}$  is a semi-group for both summation and multiplication with 0 (1 respectively) as its neutral element.  $(\mathbb{R} \cup \{NaN\}, +)$  and  $(\mathbb{R} \cup \{NaN\} \setminus \{0\}, \cdot)$  are no groups, since there is no inverse element of  $NaN$ .

Consider the example precision= $NaN$ , and recall=0, then the  $F_1$ -score becomes

$$\begin{aligned} F_1(NaN, 0) &= \frac{2NaN \cdot 0}{NaN + 0} = \frac{NaN}{NaN} = NaN \end{aligned} \tag{9}$$

Note that sometimes, the  $F_1$ -score is also defined directly based on TP, FP, and FN as follows [10]:

$$F_1(TP, FP, FN) = \frac{2TP}{2TP + FP + FN} \tag{10}$$

which leads to other behaviors regarding definedness. However, in this paper, we stick to the usual definition based on precision and recall.

#### IV. THE F-MEASURE PARADOX

Recall the four propositions from Section II.

1.  $a = 0 \Rightarrow H(a, b) = 0$
  2.  $H(a, b) = 0 \Rightarrow a = 0$
  3.  $a \neq 0, a = H(a, b) \Rightarrow b = H(a, b)$
  4.  $a \neq 0, b = H(a, b) \Rightarrow a = H(a, b)$
- (11)

If we set  $a=precision$  and  $b=recall$ , those four propositions become:

1.  $prec = 0 \Rightarrow F_1(prec, rec) = 0$
  2.  $F_1(prec, rec) = 0 \Rightarrow prec = 0$
  3.  $prec \neq 0, prec = F_1(prec, rec) \Rightarrow rec = F_1(prec, rec)$
  4.  $prec \neq 0, rec = F_1(prec, rec) \Rightarrow prec = F_1(prec, rec)$
- (12)

From Section II, one would expect that Proposition 1 and 3 are true and Proposition 2 and 4 are false. But, surprisingly, it is just the opposite. In fact, propositions 1 and 3 are false and propositions 2 and 4 are true.

*Proof:* For proposition 1, we give a counterexample. Consider the confusion matrix in Figure 2. For this matrix, the precision assumes 0 and the recall  $NaN$ . Therefore, the  $F_1$ -Score is given as  $\frac{2 \cdot 0 \cdot NaN}{2 + NaN} = NaN \neq 0$ , which concludes the proof by counterexample.

		Prediction outcome		total
		p	n	
Actual value	p'	0	8	P'
	n'	0	10000	N'
total		P	N	

Figure 2. Example confusion matrix as counterexample for proposition 1

Proposition 2: Consider the second proposition and let us assume that the  $F_1$ -Score is zero. Hence, either precision or recall is zero. In case, the precision is zero, our proof is finished. So let us, therefore, assume instead that the recall is zero. Since the  $F_1$ -score is defined (not  $NaN$ ), both recall and precision must be defined too. Furthermore, we have:

$$\begin{aligned} 0 = rec &= \frac{TP}{TP + FN} \\ \Rightarrow TP &= 0 \\ \Rightarrow \frac{TP}{TP + FP} &= 0 \\ \text{(Precision is not } NaN, \text{ therefore } TP + FP \neq 0) \\ \Rightarrow prec &= 0 \end{aligned} \tag{13}$$

■

Proposition 3: Again, we give a counterexample, we can use the same confusion matrix as for proposition 1. With this we get  $prec = NaN = F_1(prec, rec)$  and  $rec = 0$ .

Proposition 4:

*Proof:* We discern the following three cases:

Case 1:  $rec = F_1(prec, rec) = NaN$

$$\begin{aligned} rec &= F_1(prec, rec) = NaN \\ \Rightarrow TP &= 0 \\ \Rightarrow FP + TP &= 0 \\ \text{(since } prec \neq 0) \\ \Rightarrow prec &= NaN = rec = F_1(prec, rec) \end{aligned} \tag{14}$$

Case 2:  $rec = F_1(prec, rec) = 0$

Due to proposition 2, it follows that  $prec = 0 = F_1(prec, rec)$ . Since the precision cannot diminish, this case actually turns out to be impossible.

Case 3:  $rec = F_1(prec, rec) \neq 0$  and  $rec = F_1(prec, rec) \neq NaN$ .

The precision cannot assume  $NaN$ , since otherwise the  $F_1$  score would be  $NaN$ , too.

Furthermore, from the definition of the harmonic mean, it is known that:

		Prediction outcome		total
		p	n	
actual value	p'	0	0	P'
	n'	8	10000	N'
total		P	N	

Figure 3. Example confusion matrix as an counterexample for proposition 1 (case recall)

$$rec = F_1(prec, rec) = \frac{2prec \cdot rec}{prec + rec}$$

Since  $rec \neq 0$ , we can divide both sides of the equation by the recall and obtain:

$$1 = \frac{2prec}{prec + rec}$$

We then multiply both sides of the equation by  $prec + rec$ :

$$prec + rec = 2prec$$

Finally, we subtract from both sides of the equation the precision and get:

$$rec = prec \text{ which constitutes just the conducted claim. } \blacksquare$$

Finally, let us investigate if the same paradox holds also for F-measure and recall instead of precision. Analogously to the precision case, we first present a counterexample (see Figure 3).

This time, the recall is 0 and the precision NaN, which leads to a NaN  $F_1$ -Score.

*Proof:* Let us now prove the second proposition for F-measure and recall. Again, since the  $F_1$ -score is zero, neither of precision and recall can be NaN. If the recall is zero, our proof is finished. Therefore, let us instead assume the precision is zero.

$$\begin{aligned} 0 = Precision &= \frac{TP}{TP + FP} \\ \Rightarrow TP &= 0 \\ \Rightarrow \frac{TP}{TP + FN} &= 0 \\ \Rightarrow Recall &= 0 \end{aligned} \tag{15}$$

The two remaining properties 3 and 4 can be proven analogously.

### V. DISCUSSION

Purely formally seen, the computation rules for NaN values are mathematically consistent and correct and also reflect the standard procedure for computer-based  $F_1$ -Score implementations if they make use of ordinary floating-point computation logic. It remains, however, to investigate, if these rules are also

reasonable in the given context. The answer is partly yes and partly no. Consider first the case the recall is undefined (NaN), which means that the positive class never shows up in the evaluation data set. If the algorithm predicts only for a single data item the positive class, then the precision immediately turns to zero. In this case, an NaN  $F_1$ -Score seems to be the best choice.

However, if the precision is undefined (NaN), the matters look a bit different. In this case, the tested machine learning method would never predict the positive class. If the positive class shows up quite a few times in the evaluation data, such a method would clearly perform very poorly and an  $F_1$ -score of zero would seem adequate.

So far, we investigated only the  $F_1$ -score, although in the title we mentioned the F-measure in general. This generalized F-measure is given by:

$$F_\beta = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}} \tag{16}$$

For  $\beta = 2$  we get for instance the  $F_2$ -score, which is defined as:

$$F_2 = 5 \cdot \frac{\text{precision} \cdot \text{recall}}{(4 \cdot \text{precision}) + \text{recall}} \tag{17}$$

The F-measure for  $\beta \geq 2$  penalizes a poor recall stronger than a bad precision, since in some situations, like cancer detection, missing any items of the positive class can be fatal in practice. However, the use of different weighting factors does not influence in any way the properties derived here. Thus, our findings also hold for the F-measure in general.

Finally, this paradox also reveals a shortcoming of most mathematical proofs. Undefined values are not rare in practice. They can be caused by missing values or incomputability as investigated here. Albeit, in proofs, they are usually completely ignored. The paradox investigated here shows that such undefined values can easily flip statements completely around.

While the findings as stated here are mainly theoretical, they can have some practical implications as well. If the different behaviors of harmonic mean and F-measure as described here were ignored, then in certain anomalous situations, incorrect conclusions might be drawn from the data.

### VI. CONCLUSION

We presented two basic statements about the harmonic mean, where the first is true and the second false. However, for the  $F_1$ -score as the harmonic mean of precision and recall, the truth value of both statements is completely turned around. This paradox is caused by the fact that the possibility that input values can be undefined is not taken into account in the original propositions for the harmonic mean. Hence, with this paradox, we also revealed an important shortcoming of mathematical proofs in general.

### REFERENCES

- [1] S. J. Farlow, Paradoxes in Mathematics. Chicago, Illinois: Dover Books, 2014.
- [2] A. Whitehead and B. Russell, The Principles of Mathematics, 2nd ed. New York, New York: W. W. Norton & Company, 1996.

- [3] S. Banach and A. Tarski, "Sur la décomposition des ensembles de points en parties respectivement congruentes (on the decomposition of sets of points into respectively congruent parts)," *Fundamenta Mathematicae*, vol. 6, 1924, pp. 244–277.
- [4] B. Efron and C. Morris, "Stein's paradox in statistics," *Scientific American*, vol. 236, no. 5, 1977, pp. 119–127.
- [5] B. Abma, "Evaluation of requirements management tools with support for traceability-based change impact analysis," Master's Thesis, University of Twente, 2009.
- [6] J. S. Akosa, "Predictive accuracy : A misleading performance measure for highly imbalanced data," in *Proceedings of the SAS Global Forum*, 2017, p. 2–5.
- [7] D. B. MacNeil, *Fundamentals of Modern Mathematics: A Practical Review*. Dover Publications, 2013.
- [8] C. D. Manning and H. Schütze, *An Introduction to Information Retrieval*. Cambridge University Press, 2009.
- [9] L. Běhounek and M. Daňková, "Extending aggregation functions for undefined inputs," in *Proceedings of the International Symposium on Aggregation and Structures*, 2018, pp. 15–16.
- [10] D. Chicco and G. Jurman, "The advantages of the matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation," *BMC Genomics*, vol. 21, no. 6, 2020, pp. 1–13.

## Detecting Suicide Risk Through Twitter

Javier Fabra  
Aragón Institute of Engineering  
Research (I3A)  
Department of Computer Science and  
Systems Engineering  
Universidad de Zaragoza, Spain  
email: jfabra@unizar.es

Ana B. Martínez-Martínez  
IIS Aragón  
Faculty of Health Sciences  
Universidad de Zaragoza, Spain  
email: amarmar@unizar.es

Yolanda López-Del-Hoyo, María  
C. Pérez-Yus, Bárbara Oliván-  
Blázquez  
IIS Aragón  
Department of Psychology and  
Sociology  
Universidad de Zaragoza, Spain  
emails:  
{yolandal, mcperesy, bolivan}@unizar.es

**Abstract**— Mental illness is one of the main causes of illness worldwide. Currently, it is estimated that about 300 million people suffer from depression according to the World Health Organization (WHO). In this context, this work deals with the construction of a platform that allows to detect the risk of suicide using data from Twitter. This platform combines external emotional processing systems, clustering techniques and a system based on machine learning that facilitate the automatic classification of the information obtained. The entire process is articulated around a multidisciplinary team of professionals in Health Sciences and Information Technology, generating as a result a useful prototype for suicide prevention in the population.

**Keywords**— mental health; suicide; prevention; Twitter; clustering; automatic classification.

### I. INTRODUCTION

Mental illness is one of the main causes of illness worldwide. Currently, it is estimated that about 300 million people suffer from depression according to the World Health Organization (WHO). However, the provision of services for the identification, support and treatment of this type of mental illness globally is considered insufficient [1]. Although 87% of the governments of the different countries offer some type of basic care service for mental health problems, 30% of them do not have specific programs or budgets for mental health [1]. Furthermore, there are no definitive tests for the reliable diagnosis of most mental illnesses. The typical diagnosis is based on the patient's self-reported experiences, behaviors reported by family and friends, and the clinical examination of their mental state.

The data traditionally obtained through survey methodology are not a real-time reflection of the true state of mental and emotional health of individuals, which does not allow to offer a reliable estimate of the population's mental health. In Spain, suicide is the main cause of unnatural death, doubling the number of deaths in traffic accidents. The impact of suicide on families is devastating, when many of the deaths caused by suicide could be prevented [2]. Understanding how people communicate their suicidal tendencies is a cornerstone to preventing such deaths [3].

Social network platforms, such as Twitter, Instagram or Facebook are a source of faithful and real-time data on the emotional state of people [4]. In this work, we focus specifically on suicide-related aspects and propose the development of a platform capable of detecting and analyzing the emotional states of people globally and individually from the information available on social networks, specifically Twitter. The objective is twofold: on the one hand, to understand efficiently demand in the places and times that it occurs; on the other, to develop tools for suicide detection and prevention. In this first approach, we will focus on Tweets written in Spanish, although the methodology is applicable in other languages.

As the literature review points out, Twitter is one of the most widely used social media platforms worldwide, and has been the subject of numerous previous studies. Geographic, daily, weekly, and seasonal patterns of positive and negative affect have been observed in some of these studies [5][6]. The measurement of happiness levels in the populations of certain countries has also been analyzed [7]. Happiness was found to correlate with demographic and general well-being characteristics. The potential of Twitter to detect depression has also been studied [4]. The greater detection of emotional patterns specifically related to mental health variables is of special interest for global health, by helping to understand the places and moments of greatest demand (unmet) and the effective provision of resources that respond to these needs [8][9].

Previous studies have collected and classified the Tweets related to suicide [6]-[11]. However, these databases are still insufficient and the development of models for automatic detection is still rather immature. Although Twitter may provide an unprecedented opportunity to identify those at risk of suicide [10], as well as an intervention mechanism for both at the individual and community level, valid, reliable and acceptable methods for online detection have not been developed yet [12]. The best modus operandi for suicide prevention through social media remains to be clarified, so this work points to address open and existing problems in this area.

The potential of social media as data sources for improving people's health and quality of life is a relatively new phenomenon that society is beginning to value and

understand. This work represents a step within the enormous range of possibilities that the use of social media opens in the area of health. Accessing data and managing the large amount of information that can be collected (millions of daily entries) is a complex task that requires the integration of different technologies in order to integrate the various sources of information with a data processing system that allows us to obtain an adequate analysis of the data collected.

In this work, a deployment based on computer technologies for Tweet collection and analysis, as well as a system based on machine learning techniques that facilitate the automatic classification of the information obtained are depicted. This work deals with a subject and a set of technologies that have been previously considered by other researchers [13]-[16]. However, it represents a full framework, engineered and implemented using various technologies, and structured around a multidisciplinary team of professionals in Health Sciences and professionals in Information Technology. As a result, a useful prototype for suicide prevention and detection of real emotional states in the population has been developed.

The techniques developed in this work are easily adaptable to other contexts and studies in mental health and even in other sectors and institutions. The use of a technological framework based on languages and tools for the processing of information flows (streams) in real time facilitates the reuse of the main concepts and ideas underlying this work in other areas.

The remainder of this paper is as follows. Section II depicts the methodology and steps carried out and the architecture overview of our proposal. Section III presents implementation details, as well as the results from the experimentation conducted. Finally, Section IV concludes the paper and presents some related lines that currently are being addressed

## II. METHODOLOGY AND ARCHITECTURE

Figure 1 presents the methodology of the solution that we propose, and that corresponds to the workflow to be carried out.

The first step is to obtain the Twitter entries from some keywords related to suicide. To identify potentially emotional tweets, a large vocabulary of emotional terms has been compiled from different sources, including *The Spanish adaptation of Affective Norms for English words (ANEW)* and the *Spanish dictionary of the Linguistic Inquiry and Word Count (LIWC)* [17][18]. ANEW provides a set of emotional normative scales for a set of words. Furthermore, LIWC is an analysis software that calculates the degree to which people use different categories of words across a wide spectrum of texts. Validation studies reveal that LIWC satisfactorily assesses positive and negative emotions.

One of the hypotheses of our proposal is that adding properties to the text contained in the Tweet facilitates and improves the identification and classification of suicide risk groups. Therefore, a series of properties associated with the text are obtained and added below, which are based both on external natural language processing systems and platforms and on internal algorithms that obtain the information through a text evaluation platform by part of selected reviewers in the area of Health Sciences and Medicine. The emotional vocabulary has been organized by combining the hierarchy of emotions by W. G. Parrott [28] and the *tree of emotions* by Shaver et al. [29]. Each emotional word has been classified into six categories of *primary emotions* of love, joy, surprise, anger, sadness and fear, with 25 subgroups of *secondary emotions*. This task has been carried out by integrating the execution of the Indico affective and emotional text processing tool [19].

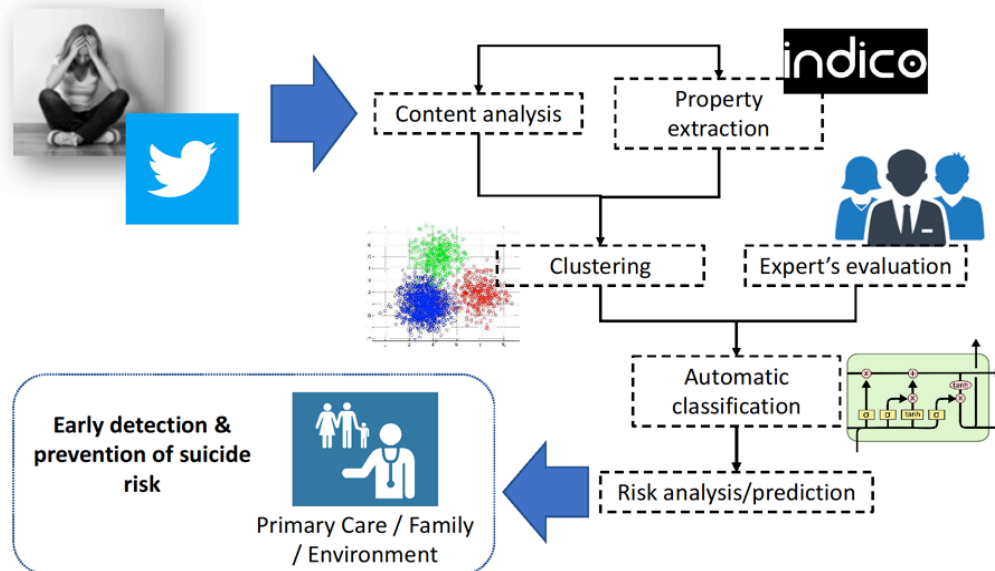


Figure 1. Methodology for early detection and prevention of suicide risk in Twitter.

Once the tweets with the properties have been obtained and annotated using Indico, a clustering is carried out to group the rich tweets. Clustering aims to generate data groups with similar characteristics. In our case, we want to identify suicide risk groups. The selected method for clustering is *k-means* [20]. This method is based on partitioning the data into *k* well-defined groups. To select this *k* value, two methods that offer good results have been used: the *elbow* method [21] and the *cross validation* [22].

The clusters that are obtained must be analyzed to know the characteristics of each one. This step must be carried out by a group of expert reviewers on the subject to make an evaluation of the clusters and to know the quality of the groups. The objective of this analysis is to validate that the clusters obtained correspond to suicide risk groups. At this stage, human coding is used to determine the degree of relationship of the classified Tweets, according to the judgment of the coding team made up of researchers from mental health and medicine. These researchers are specialized in suicide prevention and have training in detecting suicide risk.

Finally, in the last step an automatic classifier fed with the clusters and Tweets is created. This classifier is capable of receiving new Tweets and classifying them into one of the groups in order to identify whether there is a risk of suicide or not. In this work, the use of a *Long Short-Term Memory* (LSTM) neural network has been chosen. The result of this process is a tool that is capable of predicting whether a Tweet is in the suicide risk group or not. Since neural networks improve as they feed, the tool will always be in constant advance.

### III. EXPERIMENTATION AND RESULTS

In this section, we will detail the implementation and experimentation that have been carried out in a first pilot study to validate the approach presented in this work.

#### A. Obtaining and adding properties

To obtain Tweets that contain the designed vocabulary of emotional-type words described in Section II, the Amazon Web Service infrastructure has been used. Multiple instances of Elastic Compute Cloud (EC2) [23] have been used to receive the data stream through an application implemented with NodeJS and using the Twitter API [24]. The summary data has been stored in MongoDB for access from the web and API. As a result, 3,051 context-sensitive Tweets have been extracted. The Tweets with emotional information have been then annotated using the Indico API, providing more information about the characteristics of the text.

Let us to briefly describe the properties that have been extracted:

- *Positivity*: value that represents the probability that the Tweet is positive or negative, if the value is greater than or equal to 0.6 the text has a positive feeling and if it is less than 0.6 the feeling is negative.
- *Engagement*: probability that the text will be bookmarked or retweeted by other people.

- *Emotional content*: they represent five values to determine the emotions that the author has expressed in the Tweet. The five emotions that are obtained as a result are: *anger, joy, fear, sadness and surprise*.
- *Personality*: set of four values to define the author's personality traits. The four personalities are: *extroversion, sincerity, sympathy and meticulousness*.
- *People*: there are sixteen values that represent the probability that the author adjusts to one of them. People are those described by Myers Briggs [25].

#### B. Clustering and expert's evaluation

Knime [26] has been used to implement clustering, although there are other equally valid alternatives (R or Weka, among others). Knime is a platform for data analysis that also integrates components for machine learning and data mining. Developed on the Eclipse platform and programmed, mostly in Java, it has a very comfortable interface that allows you to see the work done at all times. This is possible thanks to the *workflows*. A workflow is a graphical representation in which nodes and meta-nodes, a set of encapsulated nodes, are added to read data, do operations with them, or generate output data.

The work with Knime is mostly graphic, but it also allows a high degree of configuration, addition of external code and integration with other tools, making it the most suitable candidate to carry out the clustering phase in this project.

The corresponding workflows have been implemented to both calculate the optimal *k* and to cluster using the *k-means* method. As a result of the execution of the workflows, a value of *k* = 4 was obtained as optimal.

Taking this value as optimal-*k*, the clusters were obtained. The results of the clustering process with *k-means* for *k* = 4 show a well differentiated distribution in the input collection. Table I shows some data about the Tweets that correspond to each of the clusters (C0 to C3), as well as the positivity and the five emotions detected from the contextual information of the Tweet.

TABLE I. CLUSTERS OBTAINED WITH K=4

Cluster	#Tweets	Positivity	Anger	Joy	Fear	Sadness	Surprise
#0	654	0.68	0.15	0.32	0.14	0.27	0.13
#1	884	0.80	0.24	0.23	0.15	0.26	0.12
#2	604	0.29	0.25	0.11	0.16	0.40	0.08
#3	909	0.42	0.24	0.09	0.22	0.39	0.05

As it is shown, cluster 0 groups 21% of the Tweets of the input collection (3051 Tweets), cluster 1 (C1) 29%, cluster 2 20% and cluster 3 30%. Tweets classified in both clusters 0 and 1 show a positivity above the mean (0.56), with values of 121% and 143% with respect to the mean. Clusters 2 and 3, however, group negative Tweets (0.29 and 0.42, respectively).



The analysis of the emotions contained in the context of the Tweets allows us to detail the identified clusters. Tweets in cluster 0 show a low level of *anger* compared to the other clusters (65% of the mean, 0.22, compared to 107%, 113% and 110% with respect to the mean in the other clusters). The same occurs with *joy*, where both cluster 0 and cluster 1 stand out (172% and 125% respectively compared to the mean, 0.18, and 58% and 51% of clusters 2 and 3, respectively). The *fear* analysis shows that the Tweets samples contained in cluster 3 have a value above the mean (129% above the mean, 0.17), while the other clusters remain below it (79% for cluster 0, 89% for cluster 1 and 93% for cluster 2). From the analysis of *sadness*, values are obtained in line with expectations. Cluster 2 and 3 Tweets contain values above the average, 0.33 (122% and 118%, respectively), while the values of clusters 0 and 1 are below (around 80% of the half). Finally, the *surprise* analysis shows that the Tweets of cluster 3 are the ones that show the least surprise (56% with respect to the mean, 0.1), which indicates an apathetic or lazy profile in the context. The analysis of the other properties allowed to refine the results obtained. For example, *engagement* shows similar values for the four clusters (between 95% and 103%), indicating that there is not a predominant trend in any of the groups analyzed.

Subsequently, a human coding was performed to determine the degree of relationship of the classified Tweets. Reviewers were asked to conceptualize the task as the level of concern they would have if they viewed that post on their own online social network and whether they would consider the post to require further investigation from a friend, family member, or a third party. Tweets were individually examined and coded according to a classification system validated by the research team.

The analysis of the results obtained allows us to deduce that there is a direct relationship between the identified clusters and the possible suicide risk groups. Clusters 2 and 3 have a strong emotional content that is reflected in states of anger and sadness, fear and apathy in their content. The population that is at high risk of committing suicide expresses boredom with their life, routine, very negative content (some talk of death, or of being tired of living) and fear of certain scenarios or situations (the negative sensitivity towards simple situations in the day to day increases significantly).

In conclusion, the experts considered that the identified clusters fit perfectly into a classification of potential suicide risk, establishing it as follows: cluster 0 corresponds to tweets and individuals who have a very low risk of committing suicide, cluster 1 represents a low risk, cluster 2 represents a medium risk (to be monitored, since the parameters indicate that it is closer to a high risk than a low risk), and cluster 3 a high risk. Therefore, this suggests that the subsequent phases of analysis will focus on cluster 3, since it is the indicator that this information should be analyzed with priority and work together with the Primary Care services to identify the authors of the Tweet and try to implement prevention procedures as soon as possible.

### C. Automatic classification

The LSTM neural network has been designed with Tensorflow [27]. Tensorflow is a machine learning framework developed by Google, programmed in Python and C ++. It stands out for its simplicity when it comes to building and training neural networks, but at the same time obtaining great results.

To create and train the neural network, we have developed scripts in Python. The network configuration is based on 10 LSTM hidden layers with 20 neurons in each layer. The output layer is a normal layer that produces a single output. The *loss function* measures the inconsistency between the actual values of the output and the predicted ones, in this case the mean of the absolute error is used as loss. The *optimization function* helps to minimize the loss and sight function. *Adam (Adaptive Moment Estimation)* has been used as it offers very low loss values. To evaluate the classifier, the *accuracy* will be used to obtain the percentage of correct answers in the predicted values.

Regarding training, the input data used are the properties of the Tweets together with the cluster to which they belong. These data have been divided into training data (70% of the total) and test data (30% of the total). To improve training, the data that accounted for 70% of the total has been partitioned to have 80% of that set as training data and 20% as validation data.

Each workout has been run 5 times to ensure that the result obtained is reliable. The training is done with 100 iterations, adding a field to end the training before completing them if the loss value does not improve in 3 consecutive iterations. When executing a problem appears, sometimes the best solution is not reached. To avoid this problem and ensure convergence, the 100 iterations have been left for the results shown below.

Training has been considered providing the network with all the properties in the Tweet. With  $K = 4$ , the evaluation function returns an accuracy of 93.34%. The confusion matrix obtained can be seen in Table II.

The first row represents the actual values that belong to cluster 0, the second row to cluster 1, the third row to cluster 2, and the last row to cluster 3. The columns appear in the same order. The results obtained in this training are very positive.

It can be observed that 24 data correspond to cluster 3, but belong to 1, and 16 data that belong to cluster 2 have been predicted as belonging to 3. For cluster 0, a success rate of 98.96% is obtained, for cluster 1 the percentage is 99.26%, for cluster 2 87.79% is obtained and for cluster 3 the percentage is 87.23%. In the cases of cluster 0 and cluster 1, the percentages are very close to 100%. The other two clusters do not reach 90% but they get a very high percentage as well.

The success rate that has been obtained is quite high, with an accuracy of 93.34% for the test data. Furthermore, the confusion matrix only presents a few false positives and negatives for cluster 3, assuming a fairly low percentage with respect to the successes achieved.

TABLE II. CONFUSION MATRIX

	C0	C1	C2	C3
C0	191	1	1	0
C1	1	267	1	0
C2	4	1	151	16
C3	0	24	12	246

#### IV. CONCLUSIONS

This work has addressed the problem of automatically identifying suicide risk groups in the context of Twitter. Based on a collection of Tweets obtained by keywords in relation to suicide, a series of properties have been added to enrich them. Then, clustering was applied using k-means with an optimal value of  $k = 4$ . The process has been validated by a group of experts in the context of mental health. This validation allowed establishing a relationship between the clusters obtained and the levels of suicide risk.

Finally, an automatic classifier has been built using an LSTM neural network. The neural network has been configured with 10 hidden layers and 20 neurons per layer. After training and evaluating the neural network with the test data, an accuracy of 93.34% has been obtained.

The proposal presented in this work shows very satisfactory and promising results. This approach is currently being extended, deploying the platform on an Amazon AWS infrastructure to automate the entire process and the different phases. As a result, direct connection with Primary Care Services is being worked on, so that the detection of a positive case allows initiating a series of actions to identify and contact the possible author of the content of the tweet. However, this is a very complex process that is being developed.

In addition to the work done, there are several ideas to improve the results obtained. Regarding the clustering process, other distance functions could be used. When implementing k-means the distance function that has been used is the Euclidean distance. Using other types of distance could improve the clusters obtained. A good option would be to use the Tanimoto coefficients to find the similarity and diversity of the sample set.

On the other hand, the properties could be hierarchized, studying the current properties to know which of them are more important in the generated groups. Then, we would have to add a weight to the most relevant properties so that they would have more importance when clustering or look for alternative techniques that allow us to apply these priorities.

We have also considered the possibility of using other classification techniques. Neural networks generate very good results, but there are other techniques, such as Random Forest or Support Vector Machine. It would be interesting to classify with these or other methods and compare the results between them.

#### ACKNOWLEDGMENT

This work has been supported by the JIUZ-2018-TEC-04 project, granted by Fundación Ibercaja and Universidad de Zaragoza. The authors of this paper want to specially thank David Fustero for his collaboration in the implementation of the platform; Claudia García Martínez for her help and support in conducting the study; as well as the experts in Health Sciences and Medicine for providing us with the revisions of the data used in this study, and for giving feedback on the results: Lara Barahona, Alberto Barceló, María Beltrán, Luis Borao, Roberto Buil, Daniel Campos, Luis Cortés, Irene Delgado, Paola Herrera, Laura Izquierdo, Andrea Lafuente, Andrea Llera, Marta Modrego, Alicia Monreal, Héctor Morillo, Mar Posadas, Marta Puebla, Marta Puértolas, Yaravi Rodríguez, Samara Sáez, Sara Sin, Sol Torres, David Valera, and Francisco Daniel Vinués.

#### REFERENCES

- [1] R. Detels, "Oxford textbook of Public Health", in Oxford medical publications, Oxford University Press, 2009.
- [2] J. M. Antón San Martín, "The impact of suicide on the family: the specific process of family grief", in *Redes: Revista de Psicoterapia Relacional e Intervenciones Sociales* (24), pp. 109-123, 2010.
- [3] B. O'Dea, S. Wan, P. J. Batterham, A. L. Calear, C. Paris, and H. Christensen, "Detecting suicidality on Twitter", in *Internet Interventions*, vol. 2(2), pp. 183-188, 2015.
- [4] M. De Choudhury, M. Gamon, S. Counts, and E. Horvitz, "Predicting Depression via Social Media", in the 7th International Conference on Web and Social Media (ICWSM), 2013.
- [5] P. S. Dodds, K. D. Harris, I. M. Kloumann, C. A. Bliss, and C. M. Danforth, "Temporal patterns of happiness and information", in *A Global Social Network: Hedonometrics and Twitter*. PloS one, vol. 6(12), 2011.
- [6] J. Luo, J. Du, C. Tao, H. Xu, and Y. Zhang, "Exploring temporal suicidal behavior patterns on social media: Insight from Twitter analytics", in *Health Informatics Journal*, vol. 26(2), pp. 738-752, 2020.
- [7] L. Mitchell, M. R. Frank, K. D. Harris, P. S. Dodds, and C. M. Danforth, "The Geography of Happiness: Connecting Twitter Sentiment and Expression, Demographics, and Objective Characteristics of Place", in *PLoS ONE*, vol. 8(5), pp. e64417, 2013.
- [8] M. J. Paul and M. Dredze, "Discovering Health Topics in Social Media Using Topic Models", in *PLoS ONE*, vol. 9(8), pp. e103408, 2014.
- [9] S. Alotaibi, R. Mehmood, I. Katib, O. Rana, and A. S. Albeshri, "A Big Data Analytics Tool for Healthcare Symptoms and Diseases Detection Using Twitter, Apache Spark, and Machine Learning", in *Appl. Sci.*, vol. 10, 2020.
- [10] J. Jashinsky et al., "Tracking suicide risk factors through Twitter in the US", in *Crisis*, vol. 35, pp. 51-59, 2013.
- [11] M. J. Vioulès, B. Moulahi, J. Azé, and S. Bringay, "Detection of suicide-related posts in Twitter data streams", in *IBM Journal of Research and Development*, vol. 62, no. 1, pp. 7:1-7:12, 2018.
- [12] H. Christensen, P. J. Batterham, and B. O'Dea, "E-health interventions for suicide prevention", in *International Journal of Environmental Research and Public Health*, vol. 11(8), pp. 8193-8212, 2014.

- [13] P. Burnap, G. Colombo, R. Amery, A. Hodorog, and J. Scourfield, "Multi-class machine classification of suicide-related communication on Twitter", in *Online Social Networks and Media*, vol. 2, pp. 32-44, 2017.
- [14] A. Abboute et al., "Mining Twitter for Suicide Prevention", in *Natural Language Processing and Information Systems*, vol. 8455, pp. 250-253, 2014.
- [15] S. Fodeh, J. Goulet, C. Brandt, and A. T. Hamada, "Leveraging Twitter to better identify suicide risk", in *Proceedings of The First Workshop Medical Informatics and Healthcare, 23rd SIGKDD Conference on Knowledge Discovery and Data Mining*, PMLR 69:1-7, 2017.
- [16] K. D. Varathan and N. Talib, "Suicide detection system based on Twitter", in *2014 Science and Information Conference*, pp. 785-788, 2014.
- [17] J. Redondo, I. Fraga, and I. Padrón, "The Spanish adaptation of ANEW (Affective Norms for English Words)", in *Behavior Research Methods* vol. 39, pp. 600–605, 2007.
- [18] J. Pennebaker, M. Francis, and R. Booth, "Linguistic inquiry and word count (LIWC)", 1999.
- [19] "Indico - Intelligent Process Automation for Document Intake, Understanding and Digitization", available at <https://indico.io/> [Last access: 2020-09-11]
- [20] V. Faber, "Clustering and the continuous K-means algorithm", vol. 22, Los Alamos Science, 1994.
- [21] A. Hardy, "An examination of procedures for determining the number of clusters in a data set", in *New Approaches in Classification and Data Analysis*, pp. 178-185, Springer Berlin Heidelberg, 1994.
- [22] M. Stone, "Cross-Validatory Choice and Assessment of Statistical Predictions", in *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 36(2), pp. 111-147, 1974.
- [23] "Amazon Web Services", available at <https://aws.amazon.com/es/> [Last access: 2020-09-11]
- [24] "Twitter for Developers", available at <https://developer.twitter.com/en> [Last access: 2020-09-11]
- [25] "16 Myer Briggs personalities", available at <https://www.16personalities.com/personality-types> [Last access: 2020-09-11]
- [26] "KNIME", available at <https://www.knime.com/> [Last access: 2020-09-11]
- [27] "TensorFlow - open-source framework", available at <https://www.tensorflow.org/> [Last access: 2020-09-11]
- [28] W.G. Parrott, "Emotions in Social Psychology" in *Psychology Press*, Philadelphia, 2001.
- [29] P. Shaver, J. Schwartz, D. Kirson, C. O'Connor, "Emotion knowledge: further exploration of a prototype approach", in *J. Pers. Soc. Psychol.*, vol. 52, no. 6, pp. 1061-1086, 1987.

# Adaptive Multimedia Indexing Using Naïve Bayes Classification

Clement H. C. Leung

School of Science and Engineering  
Chinese University of Hong Kong  
Shenzhen, China  
clementleung@cuhk.edu.cn

R. F. Ma & Jiayan Zhang

School of Data Science  
Chinese University of Hong Kong  
Shenzhen, China

**Abstract**— Current organizations increasingly depend on multimedia document repositories for their effective operation. However, unlike text-oriented objects, the retrieval of multimedia objects is often inhibited by limitations in their search and discovery mechanisms, since they do not readily lend themselves to automatic processing or indexing. Here, we describe the structure of an adaptive search mechanism which is able to overcome such limitations. The basic framework of the adaptive search mechanism is to capture human judgment in the course of normal usage from user queries in order to develop semantic indexes which link search terms to media objects semantics. This approach is particularly effective for the retrieval of such multimedia objects as images, sounds, and videos, where a direct analysis of the object features does not allow them to be linked to search terms, such as non-textual/icon-based search, deep semantic search, or when search terms are unknown at the time the media repository is built. An adaptive indexing mechanism is described which makes use of naïve Bayes classification approach. This approach allows for the efficient organizational creation and updating of media indexes, which is able to instill and propagate deep knowledge relating to the organizational functions into the media management system concerning the advanced search and usage of multimedia resources. The present learning approach will enable intelligent search of multimedia resources that are otherwise hard to be located and retrieved.

**Keywords** – multimedia information indexing, reinforcement learning; multi-agent; naïve Bayes classifiers; stochastic game; probability generating function.

## I. INTRODUCTION

Information search and retrieval has extended from textual based to multimedia content, with the characteristic of information search and retrieval shifting from pull to push applications. Instead of searching an accurate piece of information in a database, users are given selected choices of information [18]. In addition, affective indexing of multimedia content combines emotional responses generated by the users is sometimes employed, e.g. the psychophysiological signals, galvanic skin response, face tracking, etc, [19].

There is now general consensus that involving users in the information search and retrieval process is able to improve the overall return results [22]. In [23], it is shown that using Markov decision process improves the efficiency of locating video frames in a video, and in [24], the distribution of visual words of multimedia data is found to be probabilistic in

relation to the concept relationship formed [24]. Users often allocate the results based on some form of scoring metrics; for example, a linear combination of posterior probability is employed to refine the search results [25]. In [20], it is proposed that Reinforcement Learning (RL) approach is suitable for users exposing to raw and high-dimensional information [20], while instant rewards of the agents is generally able to impart significant improvements in the searching process [21]. In Reinforcement Learning (RL), an agent learns through the interaction with the dynamic environment to maximize its long-term rewards, in order to act optimally. Most of the time, when modeling real-world problems, the environment involved is non-stationary and noisy [1][4][6]. More precisely, the next state results from taking the same action in a specific state may not necessarily be the same but appears to be stochastic [2][7]. And the exploration strategies adopted in different categories of RL algorithms provide different levels of control to the exploration of unknown factors, which in turn give various possibilities to the learning results.

As a result, the observed rewards and punishments are often non-deterministic. For example, when one is trying to find a video for performing a particular task, a shortening of the searching time with respect to some anticipated norm may be regarded as a reward, while a lengthening of the same may be viewed as punishment. Likewise, when one is exploring a new advertising channel, a resultant significant increase in sales may be viewed as a reward, while failure to do so may be regarded as punishment. In situations like these, there are stochastic elements governing the underlying environment. In the new route to work example, whether one receives rewards or punishments depends on a variety of chance factors, such as weather condition, day of the week, and whether there happens to be road works or traffic accidents which may or may not be representative.

Noise in multimedia data is generally numerous and cannot be known or enumerated in a practical sense, and this tends to mask the underlying pattern. Indeed, if stochastic elements are absent, the learning problems involved could be greatly simplified and their presence has motivated early research in the area. As early as 1990s, mainstream research in RL, such as the influential survey assessing existing methods carried out by Kaelbling *et al.* [2], and the Explicit Explore or Exploit ( $E^3$ ) Algorithm to solve Markov Decision Process (MDP) in polynomial time [3], adopts the common assumption of a stationary environment within a RL framework. Later on, with further advances in RL, theoretical analyses addressing the concern of non-stationary environment attracted great interests. One of the works by

Brafman and Tennenholtz introduces a model-based RL algorithm R-Max to deal with stochastic games [5]. Such stochastic elements can notably increase the complexity in multi-agent systems and multi-agent tasks, where agents learn to cooperate and compete simultaneously [6][10]. Autonomous agents are required to learn new behaviors online and predict the behaviors of other agents in multi-agent systems. As other agents adapt and actively adjust their policies, the best policy for each agent would evolve dynamically, giving rise to non-stationarity [8][9].

In most of the above situations, the cost of a trial or observation to receive either a reward or punishment can be significant, and preferably, one would like to arrive at the correct conclusion by incurring minimum cost. In the case of the advertising example, the cost of advertising can be considerable and one would therefore like to minimize it while acquiring the knowledge whether such advertising channel is effective. Similarly, in RL algorithms, we are always in the hope to rapidly converge to an optimal policy with least volumes of data, calculations, learning iterations, and minimal degree of complexity [11][12]. To do so, one should explicitly define the stopping rules for specifying the conditions under which learning should terminate and a conclusion drawn as to whether the learning has been successful or not based on the observations so far.

The problem of finding termination conditions, or stopping rules, is an intensive research topic in RL, which is closely linked to the problems of optimal policies and policy convergence [13]. Traditional RL algorithms mainly aim for relatively small-scale problems with finite states and actions. The stopping rules involved are well-defined for each category of algorithms, such as utilizing Bellman Equation in  $Q$ -learning [14]. To deal with continuous action spaces or state spaces, new algorithms, such as the Cacla algorithm [15] and CMA-ES algorithm [16], are developed with specific stopping criteria. Still, most studies on stopping criteria are algorithm-oriented and do not have a unified measurement for general comparison.

In this paper, we present an approach to RL by using a naïve Bayes classification framework, which explicitly incorporates the stochastic aspects of the environment in multimedia information search and retrieval. Applying naïve Bayes methods for classification problems are often employed in a variety of contexts [26][27], such as crowdsourcing and police surveillance. Here, we shall also learn and estimate the underlying stochastic structure of the environment by making use of the random classification labels gathered in the course of the learning process. Section II presents the fundamental model of a predefined general learning policy. The information search and retrieval success based on the rewards ratio is then studied in Section III. Based on the stochastic model, Section IV analyzes the probability of exceeding cost bounds. Section V views the relative occurrences of the binary classifications from the perspective of competing multi-agents, and the final conclusions are drawn in Section VI.

## II. A PROBABILISTIC LEARNING FRAMEWORK WITH A FIXED NUMBER OF LABELS

We are concerned with a learning sequence of multimedia search and retrieval observations, each of which either results in a positive classification or negative classification. That is, we are dealing with a binary classification problem with two class labels, +1 or -1, where for convenience the former is referred to as success, and the latter, failure. Such a learning sequence in the present context corresponds to the proper association of given search terms to particular multimedia objects. We are interested in determining whether the sequential classifications indicate overall success or failure in the classification process. Evidently, if the number of +1 labels gathered is much greater than the number of -1 labels, then the conclusion drawn from the learning episode should be success, while if the opposite is true, then the corresponding conclusion should be failure. In the case of search terms to multimedia objects association, learning success would mean that the association in question is sound and should be incorporated as proper index, while failure would mean that the search term-object association cannot be established. In order to proceed with the analysis, we first let  $p$  and  $q$  (with  $p + q = 1$ ) denote the probabilities of receiving a +1 or -1 label respectively for a given classification. Furthermore, we shall make use of the naïve Bayes property that different classifications are independent of each other. Later on, we shall derive estimates for  $p$  and  $q$ , which capture the stochastic structure of the learning environment. For example, if  $p > q$ , then clearly the final conclusion should be learning success. An error often committed is that when the first few observations are all -1, one would terminate prematurely and return a verdict of failure for the learning episode. Let us consider the following learning policy; such a policy is also studied in [26, 27] and is called majority voting.

**Learning Policy I:** *On gathering a total of  $r$  labels all belonging to either +1 or -1, the learning terminates and a decision is made in accordance with the accepted margin of the majority of voting of the classifiers.*

Here, we let the random variable  $T$  represent the number of classification labeling preceding the first positive classification; i.e.  $T$  may be viewed as the waiting time to the first positive classification,

$$\Pr[T = k] = pq^k, \quad k = 0, 1, 2, 3, \dots \quad (1)$$

The probability generating function  $G(z)$  of  $T$  is given by

$$G(z) = \sum_{k=0}^{\infty} \Pr[T = k] z^k = p \sum_{k=0}^{\infty} q^k z^k = \frac{p}{(1 - qz)}. \quad (2)$$

Note that after the occurrence of the first positive classification, the process probabilistically repeats itself again, so that we have for the waiting time  $W_r$  of the  $r$ th positive classification

$$W_r = \sum_{k=1}^r T_k, \quad (3)$$

where each  $T_k$  has the same distributional characteristics as  $T$ . From [17], the probability generating function of  $G_r(z)$  corresponding to  $W_r$  may be obtained

$$G_r(z) = G_1(z)^r = \left[ \frac{p}{(1-qz)} \right]^r. \quad (4)$$

To gain a better understanding of behavior specified above, it is useful to obtain the average waiting time  $W_r$  and its variance when  $r$  positive labels are attained. From (4), the mean and variance of  $W_r$  can be derived

$$E[W_r] = G_r'(1) = \frac{rq}{p}, \quad (5)$$

$$\text{Var}[W_r] = G_r''(1) + G_r'(1) - G_r'(1)^2 = \frac{rq}{p^2}. \quad (6)$$

Furthermore, the probabilities  $\Pr[W_r = k]$  may be readily obtained from the expansion of (4) so as to study the probabilities for various waiting time,

$$\Pr[W_r = k] = \binom{-r}{k} p^r (-q)^k, \quad k = 0, 1, 2, 3, \dots \quad (7)$$

As  $W_r$  is the sum of independent identically distributed random variables, when  $r$  is appreciable, it may be approximated by the normal distribution [17]

$$W_r \sim N\left(\frac{rq}{p}, \frac{rq}{p^2}\right), \quad (8)$$

whence we have, denoting by  $\Phi$  the standard normal distribution,

$$\begin{aligned} \Pr[W_r > b] &= \int_{\frac{bp-rq}{\sqrt{rq}}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} \\ &= 1 - \Phi\left(\frac{bp-rq}{\sqrt{rq}}\right). \end{aligned} \quad (9)$$

### III. LEARNING SUCCESS BASED ON THE CLASS LABEL RATIO

Let  $\rho$  be the ratio of the average number of negative labels to the number of positive labels, we have

$$\rho(p) = \frac{E[W_r]}{r} = \frac{q}{p}. \quad (10)$$

From this, we determine the inherent stochastic structure of the environment by estimating  $p$  from actual observed labels ratio  $W/r$ , where  $W$  is the sample mean of  $W_r$ . We can form our estimator from the above just by solving for  $p$ . We shall estimate the probability  $P_b$  that the learning cost for this component exceeding this bound. From (7) above, this is given by

$$P_b = 1 - \sum_{k=0}^b \Pr[W_r = k] = 1 - \sum_{k=0}^b \binom{-r}{k} p^r (-q)^k. \quad (11)$$

Here, the normal approximation can be invoked. In many RL learning episodes,  $r$  tends to be under 100, as a lengthy iteration time is not feasible and most learning algorithms aim to converge in minimum time.

Clearly, the selection of the maximum cost weight  $b$  will have a significant impact on  $P_b$ . Very often, it is more meaningful to relate  $b$  to  $E[W_r]$  either additively or multiplicatively. Table I tabulates the values of  $P_b$  for different values of  $b$ . The first part of Table I considers  $b$  by adding a fixed value  $d$ , with  $d = 5$  and  $d = 10$ , while the second part considers  $b$  by multiplying by a fixed multiple  $\alpha$ , with  $\alpha = 1.2$  and  $\alpha = 1.5$ ; here,  $b$  is rounded to the nearest integer. In the first part of Table I, we see that for either value of  $r$ , when  $p$  is appreciably greater than  $q$ , the probability of exceeding cost bounds tends to be acceptably small, and this is especially so for  $r = 20$ . The reason is that, since  $d$  is a fixed value, its relative contribution to  $b$  increases as  $p$  increases, produces a relatively large cost bound weight compared to the average one, and accordingly lowers the probability of exceeding the bound. However, in the second part of Table I, the difference between  $E[W_r]$  and  $b$  decreases as  $E[W_r]$  decreases, so that  $P_b$

TABLE I. ANALYSIS OF PROBABILITIES OF EXCEEDING COST BOUNDS

<b>b Formula</b>	<b>r</b>	<b>p</b>	<b>q</b>	<b>E[W<sub>r</sub>]</b>	<b>b</b>	<b>P<sub>b</sub></b>	<b>P<sub>b</sub>'</b>	<b>Err</b>
<b>b = E[W<sub>r</sub>] + d (d = 5)</b>	20	0.5	0.5	20.00	25	0.215	0.186	0.029
		0.8	0.2	5.00	10	0.023	0.026	0.003
		0.9	0.1	2.22	7	0.001	0.004	0.003
	50	0.5	0.5	50.00	55	0.309	0.279	0.030
		0.8	0.2	12.50	17	0.127	0.108	0.019
		0.9	0.1	05.56	11	0.014	0.017	0.003
<b>b = E[W<sub>r</sub>] + d (d = 10)</b>	20	0.5	0.5	20.00	30	0.057	0.059	0.002
		0.8	0.2	5.00	15	0.000	0.001	0.001
		0.9	0.1	2.22	12	0.000	0.000	0.000

$b$ Formula	$r$	$p$	$q$	$E[W_r]$	$b$	$P_b$	$P_b'$	Err
	50	0.5	0.5	50.00	60	0.159	0.147	0.012
		0.8	0.2	12.50	22	0.008	0.011	0.003
		0.9	0.1	05.56	16	0.000	0.000	0.000
$b = \alpha E[W_r]$ ( $\alpha = 1.2$ )	20	0.5	0.5	20.00	24	0.264	0.226	0.038
		0.8	0.2	5.00	6	0.345	0.253	0.092
		0.9	0.1	2.22	2	0.556	0.380	0.176
	50	0.5	0.5	50.00	50	0.159	0.147	0.012
		0.8	0.2	12.50	15	0.264	0.215	0.049
		0.9	0.1	05.56	7	0.280	0.207	0.073
$b = \alpha E[W_r]$ ( $\alpha = 1.5$ )	20	0.5	0.5	20.00	30	0.057	0.059	0.002
		0.8	0.2	5.00	7	0.212	0.156	0.056
		0.9	0.1	2.22	3	0.310	0.193	0.117
	50	0.5	0.5	50.00	75	0.006	0.010	0.004
		0.8	0.2	12.50	19	0.050	0.048	0.002
		0.9	0.1	05.56	8	0.163	0.121	0.042

tends to be large for higher values of  $p$ .

In Table I, column  $P_b'$  gives the exact calculation using (11), while column  $P_b$  employs the normal approximation using (9). The absolute error between the exact calculation and the normal approximation is given by column *Err*. We see that the normal approximation is quite acceptable in most cases with absolute error less than 0.1. Note that no matter whether having  $b$  additively or multiplicatively related to  $E[W_r]$ , a higher value of  $d$  or  $\alpha$  always gives smaller absolute error. We therefore suggest that the approximation should only be used when  $r$ ,  $d$  and  $\alpha$  are sufficiently large.

#### IV. MULTI-AGENT LEARNING

In *Learning Policy I* above, the termination of a learning episode is triggered whenever a fixed number of positive labels  $r$  is obtained, irrespective of the number of negative labels accumulated in the process of doing so. Sometimes, however, this may not be desirable, especially when an inordinate number of negative labels have been accumulated, in which case, termination should take place earlier along with the conclusion of learning failure. Therefore, one is comparing the number of positive labels gathered against the number of negative labels, and the learning is concluded as success or failure according to which of these achieve the majority.

More precisely, this may be viewed as a multi-agent tournament with two competing agents  $A$  and  $B$ , in which  $A$  is responsible for giving out the positive labels, while  $B$ , the negative labels. This framework is not unlike the game theoretic approach in statistical decision theory, where both the statistician and nature are regarded as players in the game of estimation, and also this may be regarded as a kind of stochastic game [5]. While we shall focus on the agents  $A$  and  $B$ , we note that there is a further agent, the learner, so that three agents exist in this situation. Here, when a classification results in a positive labels, then  $A$  would gain a score of one, while when an observation results in a negative labels, then  $B$

would gain a score of one. When either  $\pm 1$  label first reaches a given threshold  $h$ , then this will trigger a termination and the learning episode is concluded as success or failure according to which agent attains the threshold score first. Therefore, we have the following stopping rule:

**Learning Policy II:** *The learning process terminates when either agent, A or B, first reach the threshold of accumulating  $h + 1$  or  $-1$  classifications, which can be concluded as a success or a failure according to which agent attains the threshold first.*

Here, without loss of generality, we shall let  $h = 2m+1$  be odd, where  $m$  is an integer, and similar to Section II, we let  $p$  and  $q$ , with  $p + q = 1$ , signify the probabilities of receiving a positive labels, and negative labels, respectively for a particular classification. In other words, for a given classification, agent  $A$  wins with probability  $p$ , while agent  $B$  wins with probability  $q$ . In order to attain  $h$  for either agent, the number of classifications  $\Omega$  will fall within the range

$$2m + 1 \leq \Omega \leq 4m + 1 .$$

If  $f_k$  represents the probability that  $A$  wins at classifications number  $4m+1-k$ , which occurs if and only if  $A$  scored  $2m$  successes in the first  $4m-k$  observations, and subsequently score a final success, then  $f_k$  is given by

$$f_k = \binom{4m - k}{2m} p^{2m+1} q^{2m-k} .$$

The probability that  $A$  reaches the threshold first, irrespective of the classification number, is therefore given by

$$P_m = \sum_{k=0}^{2m} f_k = \sum_{k=0}^{2m} \binom{4m - k}{2m} p^{2m+1} q^{2m-k} .$$

That is,  $P_m$  gives the probability that the learning is successful (i.e. agent  $A$  wins) according to *Rule B*.

Table II computes  $P_m$  for different values of  $p$ ,  $q$ , and  $m$ . We see that, as expected, when  $p = q = 1/2$ ,  $P_m = 1/2$ , since neither  $A$  nor  $B$  has any advantage over its opponent. As  $p$  increases, however,  $P_m$  will increase, reaching almost certainty as  $p$  increases beyond 0.8. If we regard  $p$  as a

TABLE II. PROBABILITIES OF LEARNING SUCCESS

$m$	$p$	$q$	$P_m$	$m$	$p$	$q$	$P_m$
1	0.5	0.5	0.5000	5	0.5	0.5	0.5000
	0.6	0.4	0.6826		0.6	0.4	0.8256
	0.7	0.3	0.8369		0.7	0.3	0.9736
	0.8	0.2	0.9421		0.8	0.2	0.9990
	0.9	0.1	0.9914		0.9	0.1	1.0000
2	0.5	0.5	0.5000	10	0.5	0.5	0.5000
	0.6	0.4	0.7334		0.6	0.4	0.9035

	0.7	0.3	0.9012		0.7	0.3	0.9964
	0.8	0.2	0.9804		0.8	0.2	1.0000
	0.9	0.1	0.9991		0.9	0.1	1.0000

measure of  $A$ 's winning ability per trial, then when  $p \gg q$ , most trials will be scored by  $A$ , so that winning the entire game (i.e. reaching  $h$  first) is almost a certainty, and this is especially so for higher values of  $h$ . It is interesting to see that when  $h$  or  $m$  is sufficiently high (e.g.  $m=10$ ), a moderate advantage for  $A$  (e.g.  $p = 0.6$ ) is enough to almost guarantee success. On the other hand,  $1-P_m$  gives the probability that agent  $B$  wins, where the measure of  $B$ 's winning probability per trial is given by  $q$ . For instance, when  $q=0.4$ , then  $B$  stands a chance of around 27% of winning the game when  $m=2$ , and a chance of winning of around 10% when  $m=10$ .

Returning to the estimation problem, by observing  $P_m$ , i.e. by computing the observed proportion of time that agent  $A$  wins, it is possible to infer the underlying probability  $p$ . While unlike in Section II, where an explicit formula exists linking directly the observations to the estimate, such explicit relationship is not available here. Nevertheless, as can be observed from Table II, useful estimation bounds can be drawn to determine whether  $p > 1/2$  or  $p < 1/2$ . We see that it is quite reasonable to estimate  $\hat{p} > 1/2$  whenever  $P_m > 1/2$ , and this would seem sufficient for most practical purposes.

## V. CONCLUSION

Since multimedia information search environments are often noisy and seldom static nor deterministic, the use of stochastic methods is therefore an unavoidable necessity. Indeed, if stochastic elements are absent, the same outcome will always occur, obviating the need for repeated observations.

In this paper, we first consider a situation where the cumulative number of classifications is pre-specified and fixed, which constitute the criterion for stopping the learning process. By observing the random positive to negative labels ratio, a meaningful estimation of either learning success or failure may be arrived at. In most practical situations, the cost of securing a classification can be significant, and this has been incorporated into our model, with the probabilities of exceeding the classifications cost bounds also derived.

We also consider a multi-agent framework where the handing out of positive and negative labels are viewed as being performed by agents. Thus, the final learning outcome is determined by a kind of stochastic game with the agents competing against each other. The termination criterion here is determined by when and how the game is won. The respective probabilities of learning success and failure are also explicitly derived. Closed-form expressions of other relevant measures of interest are obtained. A procedure for estimating the underlying stochastic structure from the observed random agent winning frequencies is also employed.

In this study, we have adopted the naïve Bayes assumption and assumed that positive labels and negative labels occur

independently. In future, it may be useful to relax this assumption and incorporate single-step or multi-step Markov dependency into the analysis. It is likely, however, that the corresponding estimation procedures will be considerably more involved.

## REFERENCES

- [1] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey, "Maximum Entropy Inverse Reinforcement Learning," *Proc. Twenty-Third AAAI Conference on Artificial Intelligence (AAAI 08)*, vol. 8, pp. 1433-1438, 2008.
- [2] L. P. Kaelbling, M. L. Littman, and A. W. Moore, "Reinforcement learning: A survey" *Journal of artificial intelligence research*, vol. 4, pp. 237-285, 1996.
- [3] M. Kearns and S. Singh, "Near-optimal reinforcement learning in polynomial time," *In Int. Conf. on Machine Learning*, 1998.
- [4] H. Santana, G. Ramalho, V. Corruble, and B. Ratitch, "Multi-agent patrolling with reinforcement learning," *Proc. Third International Joint Conference on Autonomous Agents and Multiagent Systems*, vol. 3, pp. 1122-1129, IEEE Computer Society, 2004.
- [5] R. I. Brafman and M. Tennenholtz, "R-max-a general polynomial time algorithm for near-optimal reinforcement learning," *Journal of Machine Learning Research*, vol. 3, pp.213-231, 2002.
- [6] L. Panait and S. Luke, "Cooperative multi-agent learning: The state of the art," *Autonomous agents and multi-agent systems*, vol. 11, no. 3, pp. 387-434, 2005.
- [7] E. Ipek, O. Mutlu, J. F. Martínez, and R. Caruana, "Self-optimizing memory controllers: A reinforcement learning approach," *ACM SIGARCH Computer Architecture News*, vol. 36, no. 3, IEEE Computer Society, 2008.
- [8] L. Busoni, R. Babuska, and B. De Schutter, "A comprehensive survey of multiagent reinforcement learning," *IEEE Transactions on Systems, Man, And Cybernetics-Part C: Applications and Reviews*, vol. 38, no. 2 2, 2008.
- [9] S. V. Albrecht, and P. Stone, "Autonomous agents modelling other agents: A comprehensive survey and open problems," *Artificial Intelligence* 258, pp. 66-95, 2018.
- [10] A. Tampuu, T. Matiisen, D. Kodelja, I. Kuzovkin, K. Korjus, J. Aru, and R. Vicente, "Multiagent cooperation and competition with deep reinforcement learning," *PloS one*, vol. 12, no. 4: e0172395, 2017.
- [11] A.W. Moore and C.G. Atkeson, "Prioritized sweeping: Reinforcement learning with less data and less time," *Machine learning*, vol. 13, no.1, pp. 103-130, 1993.
- [12] E. Brochu, V. M. Cora, and N. De Freitas, "A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning," unpublished.
- [13] Q. Wei, F. L. Lewis, Q. Sun, P. Yan, and R. Song, "Discrete-time deterministic Q-learning: A novel convergence analysis," *IEEE transactions on cybernetics*, vol. 47, no. 5, pp. 1224-1237, 2017.
- [14] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning* 8.3-4 pp. 279-292, 1992.
- [15] H. Van Hasselt and M.A. Wiering, "Using continuous action spaces to solve discrete problems," *Proc. International Joint Conference on Neural Networks (IJCNN 09)*, pp. 1149-1156. IEEE, 2009.
- [16] N. Hansen, S. D. Müller, and P. Koumoutsakos, "Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (CMA-ES)," *Evolutionary computation*, vol. 11, no. 1 pp. 1-18, 2003.
- [17] W. Feller, *An Introduction to Probability Theory and its Applications*, vol. 1, 3rd Edition, Wiley & Sons, 1968.
- [18] Q. Huang, A. Puri, Z. Liu, "Multimedia search and retrieval: new concepts, system implementation, and application". *IEEE transactions on circuits and systems for video technology*, 2000, 10.5: 679-692.



- [19] R.Gupta, M. Khomami Abadi, J. A.Cárdenes Cabré, F. Morreale, T. H. Falk, and N. Sebe, "A quality adaptive multimodal affect recognition system for user-centric multimedia indexing". In: Proceedings of the 2016 ACM on international conference on multimedia retrieval. ACM, p. 317-320, 2016.
- [20] Arulkumaran, K., Deisenroth, M. P., Brundage, M., & Bharath, A. A., "Deep reinforcement learning: A brief survey," *IEEE Signal Processing Magazine*, 34(6), 26-38, 2017.
- [21] X. Yao, J. Du, N. Zhou, and C. Chen, "Microblog Search Based on Deep Reinforcement Learning," *In Proceedings of 2018 Chinese Intelligent Systems Conference* (pp. 23-32). Springer, Singapore, 2019.
- [22] Y.C. Wu, T. H.Lin, Y. D. Chen, H. Y Lee, and L. S. Lee, "Interactive spoken content retrieval by deep reinforcement learning". arXiv preprint arXiv:1609.05234, 2016.
- [23] S. Lan, R. Panda, Q. Zhu, and A. K. Roy-Chowdhury, "FFNet: Video fast-forwarding via reinforcement learning", *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6771-6780, 2018.
- [24] R. Hong, Y. Yang, M. Wang, and X. S. Hua, "Learning visual semantic relationships for efficient visual retrieval", *IEEE Transactions on Big Data*, 1(4), pp.152-161, 2015.
- [25] R. Yan, A. Hauptmann, and R. Jin, "Multimedia search with pseudo-relevance feedback. In International Conference on Image and Video Retrieval" (pp. 238-247). Springer, Berlin, Heidelberg, 2003.
- [26] E. Manino, L. Tran-Thanh, and N. R. Jennings. On the Efficiency of Data Collection for Multiple Naïve Bayes Classifiers. *Artificial Intelligence*, 275: 356–378, 2019.
- [27] E. Manino, L. Tran-Thanh, and N. R. Jennings. On the Efficiency of Data Collection for Crowdsourced Classification. In Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, pp. 1568-1575, 2018.