

Automatically Dismantling Online Dating Fraud*

Guillermo Suarez-Tangil^{*}, Matthew Edwards[‡], Claudia Peersman[‡],
Gianluca Stringhini[†], Awais Rashid[‡], and Monica Whitty⁺

^{*}King’s College London, [‡]University of Bristol

[†]Boston University, ⁺University of Melbourne

Abstract

Online romance scams are a prevalent form of mass-marketing fraud in the West, and yet few studies have addressed the technical or data-driven responses to this problem. In this type of scam, fraudsters craft fake profiles and manually interact with their victims. Because of the characteristics of this type of fraud and of how dating sites operate, traditional detection methods (e.g., those used in spam filtering) are ineffective. In this paper, we present the results of a multi-pronged investigation into the archetype of online dating profiles used in this form of fraud, including their use of demographics, profile descriptions, and images, shedding light on both the strategies deployed by scammers to appeal to victims and the traits of victims themselves. Further, in response to the severe financial and psychological harm caused by dating fraud, we develop a system to detect romance scammers on online dating platforms.

Our work presents the first system for automatically detecting this fraud. Our aim is to provide an early detection system to stop romance scammers as they create fraudulent profiles or before they engage with potential victims. Previous research has indicated that the victims of romance scams score highly on scales for idealized romantic beliefs. We combine a range of structured, unstructured, and deep-learned features that capture these beliefs. No prior work has fully analyzed whether these notions of romance introduce traits that could be leveraged to build a detection system. Our ensemble machine-learning approach is robust to the omission of profile details and performs at high accuracy (97%). The system enables development of automated tools for dating site providers and individual users.

1 Introduction

The online romance scam is a prevalent form of mass-marketing fraud in many Western countries [22, 24, 8, 32]. Cybercriminals set up a false user profile on dating websites or similar online platforms (e.g., social networking sites, instant messaging platforms) to contact potential victims, posing as attractive and desirable partners [31]. Once contact has been established, scammers apply a range of techniques to exploit their victims. In many cases, they engage in a long-term fictitious romantic relationship to gain their victims’ trust and to repeatedly defraud them of large sums of money [9].

Recently, the FBI [26] reported a total loss of \$85 million through online romance scams in the US. On an individual level, IC3 complaint data showed that an average of \$14,000 was lost per reported incident of online dating fraud. Furthermore, many victims find it difficult to seek support due to being left traumatized by the loss of the relationship, and suffer from the stigma of being an online dating fraud victim [32].

Despite the magnitude of this type of cybercrime, there is an absence of academic literature on the practical methods for detecting romance scammers. Previous work has mentioned that online dating sites are employing both automated and manual mechanisms to detect fake accounts, but do not discuss the specifics [9, 4]. Some dating sites are known to use static information such as blacklists of IP addresses or proxies to identify alleged scammers [23]. However, these countermeasures can easily be evaded through e.g., low-cost proxy services using compromised hosts in residential address spaces. Dealing with online dating fraud is challenging, mainly because such scams are not usually run in large-scale campaigns, nor are they generated automatically. As a result, they cannot be identified by the similarity-detection methods used for spam filtering. Dating websites are designed for connecting strangers and meeting new people, which renders the concept of *unsolicited* messages—a key element of most state-of-the-art anti-spam systems—strategically useless [9]. Finally, romance scammers will send a series of ordinary, personalized communications to gain their victims’ trust. These communications highly resemble messages between genuine dating site users. In many cases, the actual scam is performed after a few weeks or months and after communication has moved to other, unmonitored media [31]. Therefore, it is essential to identify romance scammers before they strike.

Given that the online dating profile is the launching point for the scam, it is important to learn a) how scammers craft profiles to draw in potential victims and b) if there are any distinguishing features of these profiles which can be identified for automatic detection. This is a distinct problem from the detection of Sybil attacks or cloned profiles [19], existing methods for which typically rely upon graph-based defences or markers of automated behaviour, neither of which are applicable here. Previous research has indicated that the victims of romance scams score highly on scales for idealized romantic beliefs [2]. Thus, a scammer profile might be expected to exploit these notions of romance when designing their dating profile. To the best of our knowledge, no prior work has analyzed how these notions of romance appear as traits in dating profiles, or whether these traits can be leveraged to build a scammer detection system.

*A shorter version of this paper appears in IEEE Transactions on Information Forensics and Security. This is the full version.

In this paper, we present a machine-learning solution that addresses the detection of online dating fraud in a fully automated fashion, which is widely applicable across the dating site market—including by the users themselves. More specifically, we combine advanced text categorization and image analysis techniques to extract useful information from a large dataset of online dating user profiles and to automatically identify scammer profiles. The key contributions of our work are as follows:

- We leverage a large public database of romance scammer profiles, in combination with a large random sample of public profiles from a matched online dating site to understand the characteristic distinctions between scammer profiles and those of regular users.
- We design three independent classification modules which analyze different aspects of public profile characteristics.
- We synthesize the individual classifiers into a highly accurate ensemble classification system. Even when parts of the profile information are omitted, our system can reliably distinguish between scammer and real user profiles (F1= 94.5%, ACC= 97%), resulting in a solid solution which should generalise well to other dating sites.

To enable replication and foster research we make our tool publicly available at <https://github.com/gsuareztangil/automatic-romancescam-digger>. This paper proceeds as follows. In §2 we describe our dataset and the observed characteristics of real and scammer dating profiles. In §3 we discuss the architecture of our ensemble classification system. In §4 we detail the division of the data for training, test and validation purposes, and present our results, before discussing related work in §6 and concluding with final remarks in §7.

2 Characterizing Dating Profiles

Though variations exist within the market, typical dating profiles consist of at least one image of the user, some basic information about their key attributes, and a self-description used as a ‘sales pitch’. Our approach to scam detection focuses on these common profile components, which are present across the market. In what follows, we compare the characteristics of real dating profiles with those which are designed by scammers, and detail what we can learn about scammer targets and strategies from the differences between them.

2.1 Data

The data we use comes from a dating site `datingnmore.com`, and the connected public scamlist at `scamdigger.com`. This dating site distinguishes itself from the market on the basis of its lack of romance scammers. It screens its registrants and members to identify scams, which are then listed openly to warn the general public and anyone whose likeness may be being appropriated by the scammers. Our dataset combines an exhaustive scrape of the scammer profiles and a large random sample of one-third of the ordinary dating profiles as of March, 2017. In total, our dataset is composed of 14,720 ordinary profiles, and

5,402 scammer profiles¹. The sampling of the dating site was spread over a member index sorted by registration date, to ensure comparison with the scamlist compiled over the site’s operation.

All data used in this paper is publicly available, with no requirement to register, log in or deceptively interact with users of the dating site to collect it. Nevertheless, in the interests of privacy, no personally identifying information is revealed in this paper, including that of reported scammers. To enable replication of our results, we make available two scripts which implement the data-harvesting process that created our dataset. This enables replication while allowing dating site users to “withdraw” from future study by removing their profile from public view. The research was approved by the relevant Institutional Review Board (IRB).

The attributes available for scammer profiles and genuine profiles are slightly different. The scamlist profiles include the IP address, email address and phone number used by the scammer when registering, along with bookkeeping information on the justification used for the decision that a profile is a scam. These variables are not present in the public member information. Contrarily, there are attributes visible on public member pages—related to the dating interests of members—which were not duplicated to the scamlist for scammer profiles. For the purposes of informing discriminative and widely-applicable classifiers, we focus on those attributes which are available for both types of profile, which we divide into the following three groups:

- **Demographics:** Simple categorical information relating to the user, such as age, gender, ethnicity, etc.
- **Images:** One or more images of the user. The dating site mandates that only images showing your own face may be used as an avatar, and users are usually motivated to include pictures that illustrate their hobbies.
- **Description:** A short textual self-description from the user, in which they advertise their key traits and interests.

Different techniques are required to extract meaningful information from these profile attributes. In the following, we cover the preprocessing required for each group, and the notable features of scammer and real dating profiles.

2.2 Profile Demographics

Dating site demographics act as a filter for users. At the crudest level, most users will be searching for a particular gender of partner. Typically, age and other information about a person will also play a role in their match candidacy. In response to such filtering, users may withhold or lie about certain demographic characteristics to make themselves seem more desirable to potential partners. For most real dating site users, any such deceptions or omissions must be low-level, as they intend for a personal relationship to result [7]. Romance scammers, however, have no expectations of a real relationship, and are highly motivated to

¹There were roughly 3,500 scammer profiles in the original data, but these included ‘or’ values where specific attributes which annotators had seen the profile present differently were given multiple values. We exploded these ‘or’ attributes into different profiles to analyze the profile-variants, but in all analysis that follows were careful to avoid assigning profile-variants from one original profile to different sets or folds.

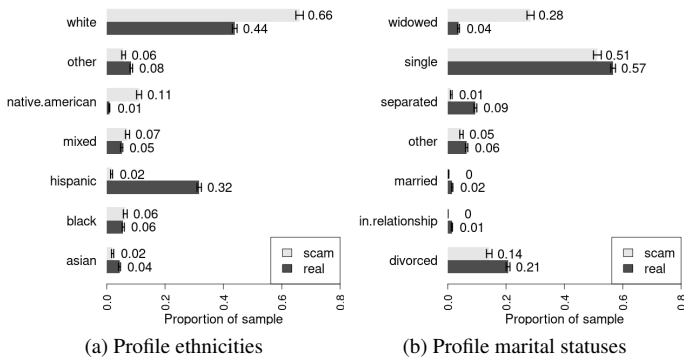


Figure 1: Ethnicity and marital status, with 95% CIs

engage in this form of deception. The information they present in profiles should thus in no way be taken as an accurate measure of their true demographics. However, the attributes selected in their profiles can reveal much about their overall strategies for attracting potential victims for romance fraud, and even, implicitly, who their targets may be.

Age, Gender, Ethnicity and Marital Status The gender distribution of both real and scam profiles was identical: about 60% of profiles are male. This highlights that romance scamming is not a gender-specific problem, in line with the understanding of previous studies [31]. The average age of real and scam profiles was around 40 in both cases. However, the distribution of ages differs significantly. Within real profiles, the average age of male and female profiles is the same, but within scam profiles the average age of females is roughly 30, and the age of male profiles is roughly 50. This bimodal distribution around the mean of real profile ages points at scammer understanding of gendered dating preferences—men here prefer younger, physically attractive partners, while women prefer partners with higher socio-economic status, who may be older [11].

As reported in Fig. 1a, the ethnicities claimed by scammers are intriguing. The high proportion which claim to be white is unsurprising, as this is the ethnicity of most of their intended victims. However, the dating site has a large Hispanic population, but the scammers rarely pretend to be Hispanic. Instead, the second most popular ethnicity amongst scam profiles is Native American, a very small population amongst the real data. This would seem to reflect some criteria of desirability which perhaps is related to the fact that dating scams are often targeted at a US-based population.

As Fig. 1b displays, while both real and scam users were mostly single, scammers prefer to present themselves as widowed rather than any of the other categories. This is unsurprising, as female scam victims often talk about such a trait being a successful strategy to gain their sympathy and trust [31]. Less desirable statuses such as divorced or separated were underrepresented in scam profiles, and scammers were far less likely than real users to be married or in a relationship.

Occupation There were a wide variety of occupations, several being misspellings or rephrases of others. Responses were grouped into 45 occupation areas. Tables 1a and 1b reflect the

TABLE I: Topmost occupation areas by presented gender

| (a) Male profiles | | | | (b) Female profiles | | | |
|-------------------|------|----------|------|---------------------|------|----------|------|
| Real | Freq | Scam | Freq | Real | Freq | Scam | Freq |
| other | 0.15 | military | 0.25 | other | 0.15 | student | 0.21 |
| self | 0.07 | engineer | 0.25 | student | 0.10 | self | 0.16 |
| engineer | 0.07 | self | 0.10 | carer | 0.08 | carer | 0.10 |
| tech. | 0.05 | business | 0.06 | service | 0.06 | sales | 0.07 |
| student | 0.05 | building | 0.06 | clerical | 0.06 | military | 0.05 |
| retired | 0.05 | other | 0.04 | teacher | 0.06 | fashion | 0.04 |
| building | 0.05 | contract | 0.04 | retired | 0.05 | business | 0.04 |
| service | 0.04 | medical | 0.03 | self | 0.04 | other | 0.04 |
| transport | 0.04 | manager | 0.02 | medical | 0.04 | finance | 0.03 |
| manual | 0.03 | sales | 0.02 | housewife | 0.03 | service | 0.03 |

major occupation areas for male and female profiles respectively. In both cases, approx. 15% of real and 4% of scam responses were not well-captured by occupation groupings, this category of ‘other’ reflecting a long tail of unique occupation responses. For both males and females, the most frequent occupation response for real profiles was “retired”, a value which was extremely rare in scam profiles.

Table 1a presents a strong bias of scam profiles towards military and engineering professions. The desirability of male military profiles is a bias romance scammers are already well-known for exploiting [31]. The masculine and high-status image of engineering might similarly explain its use by scammers. Other professions listed display a similar approach: business (in many cases, the raw response being “businessman”), medicine (i.e., “doctor”) and contracting professions which might lend themselves to explanations for why a person would later require money to be sent overseas. As shown in Table 1b, female scam profiles present less clearly suspicious occupations, with ‘student’ and ‘carer’ groups leading. The appearance of ‘fashion’ further down the list does speak towards a desirability bias (e.g., “model”). The ‘military’ group makes a surprising appearance—no real female profiles claimed such a role—even more oddly, this occupation is selected mostly by female profiles aged over 40. This may be an attempt to generalize the “military scam” used in male profiles, but its strategy is unclear. For the most part, female scam occupations fit with previous suggestions that scammers are exploiting the desirability of a young, dependent female partner in low-paying or non-professional work [31]. This role naturally lends itself to an explanation for why a person might need financial support.

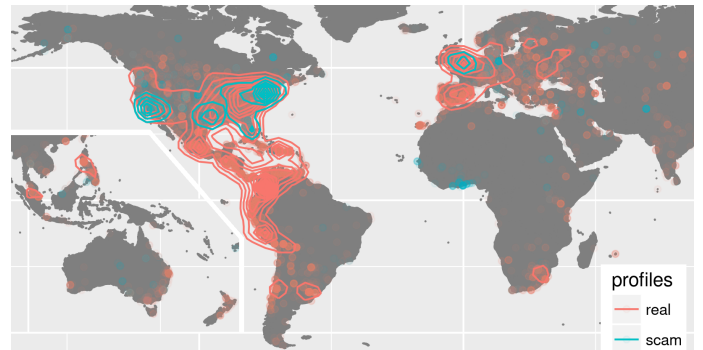


Figure 2: Worldwide location of scammer and real profiles



Figure 3: Certain contexts such as the military, academic or the medical one are often used to attract vulnerable users².

Location The location data reported in profiles was usually given to the city level, although the specificity did vary, particularly within scammer profiles. The original location responses were geocoded to provide lat/lon points, and country of origin factors. As shown in Fig. 2, the scam profiles mostly claim to be in the US or Western Europe. Corresponding with the earlier observation of a low incidence of claimed Hispanic ethnicity, scammer profiles rarely claim to be located in Latin America or Spain, despite a large real user population from these areas. This suggests that a substantial Spanish-speaking population of the dating site is not (yet) being targeted, possibly due to language barriers.

With regard to the targeted national locations, the concentration of scammer profiles in the US is highly notable, with nearly three-quarters of scam profiles with given locations claiming to be resident there. The secondary targets were the UK and Germany. More plausibly honest responses, such as Ghana, may be reactions to the dating site’s methodology of comparing IP geolocations to declared location [4]. The distribution suggests scammers are targeting rich, Western and mostly English-speaking nations.

Scammer profiles most often declared locations which were well-known Western cities. The most frequent city response was New York, being roughly 13% of all scammer locations, followed by Los Angeles (7%) and then London, Dallas, Miami, Houston and Berlin. Selecting well-known and large cities avoids the need for intricate knowledge of a smaller city/town and makes it easier for a scammer to remotely obtain enough detail to appear plausible. This approach also enables a travel narrative—commonly, a wealthy businessman originally from a large city but currently away on business.

2.3 Image Recognition

The use of images plays an important role in online dating sites. The right set of pictures can both maximize the number of people interested in a profile and help users to limit interactions to a certain type of person. This can be leveraged by criminals to reach a larger number of potential victims and to attract vulnerable users.

Scammers typically select the face of a publicly available im-

²Publicly available images from <https://www.romancescam.com/>



Figure 4: Sample of images from our dataset. Faces from *real* profiles have been redacted to preserve the anonymity.


age and build a fake persona with other images from different desirable contexts. The military context is recurrently exploited [31] as certain vulnerable victims seek a ‘knight in shining armor’. Other popular contexts are the academic and medical ones, where scammers pretend to be practitioners, students or patients. Fig. 3 shows how a scammer faked an image of a high-ranked military officer to impersonate a third party. We next study how to extract semantics from images to better understand the choices made by both legitimate users and scammers when selecting profile pictures.

As mentioned earlier, the data available in our dataset for the *scammers* category and the *real* category differs slightly. For the scammers dataset there are samples available with multiple images per profile. Conversely, for the real dataset there is usually only one picture per profile. Overall, we found images in approximately 65% of the profiles³. While the proportion of profiles with images is equivalent in both types of categories, the absolute number of images per profile is larger in the scammers dataset. Specifically, there are 0.65 images per profile in the real dataset and 1.5 in the scammers one. Note that there may be variation in the distribution of images per profile category across different dating sites. However, in our dataset, fraudsters tend to share more information than legitimate users.

A common theme found in both types of profiles is the use of pictures where users not only show their physical appearance, but also convey a sense of the hobbies or interests the person holds. Fig. 4 shows four examples of images found in our dataset where subjects are, for instance, riding or sailing.

As it is relevant how scammers present themselves in their profile pictures, we next elaborate on how to extract meaningful information from them. Recent work in the field of computer vision [10, 28] has shown that it is possible to automatically describe the content of an image accurately using deep learning. The key idea is to train a deep network with a large corpus of images for which there is a ground truth of visible context. The resulting network is then expected to i) know how to recognize elements appearing in the images and ii) be capable of generating an adequate description.

For the purpose of this work, we mainly rely on [28] to build a generative model based on a deep Neural Network (NN). The system consists of a convolutional NN combined with a language-generating recurrent NN. The model has been built using a very extensive dataset distributed by Microsoft called COCO (Common Objects in Context)⁴ with over 300,000 images. The output of the system is a meaningful description of the

³Without including the default site avatar: .

⁴<http://mscoco.org/>

image given as input.

For each image in a profile, we output the description that best represents (according to the model) the semantics involved in the picture. Fig. 4 shows images from four different profiles, two in the *real* category and two in the *scammer* category. We next show the full output extracted from the images shown in Fig. 4. For each image we output three possible descriptions with probability p .

| Descriptions automatically generated from Fig. 4a: | |
|--|--------------------|
| 1. A man riding a motorcycle down a street | $(p = 72.2e - 04)$ |
| 2. A man riding a bike down a street | $(p = 29.3e - 04)$ |
| 3. A man riding a bike down the street | $(p = 3.7e - 04)$ |

The descriptions shown above have been extracted from the image of a profile belonging to the *real* category. The image shows a man standing over a bicycle in the street. The description in 1) guessed that the man is riding a motorcycle. This misconception can most likely be attributed to the headlight and the ad banner over it (uncommon in bikes). Descriptions 2) and 3) however are guessed correctly with a probability of the same magnitude. We argue that this type of mistake is orthogonal to our problem. Confusing objects of similar kinds should not have a negative impact as long as the main activity is correctly inferred (i.e.: a man riding down the street).

| Descriptions automatically generated from Fig. 4b: | |
|--|--------------------|
| 1. A man standing in a boat in the water | $(p = 28.0e - 05)$ |
| 2. A man standing in a boat in a body of water | $(p = 9.9e - 05)$ |
| 3. A man in a suit and tie standing in the water | $(p = 2.1e - 05)$ |

The afore set of descriptions also belong to the *real* category. The image shows a man standing in the deck of a boat as correctly predicted.

| Descriptions automatically generated from Fig. 4c: | |
|--|--------------------|
| 1. A man riding on the back of a brown horse | $(p = 11.8e - 03)$ |
| 2. A man riding on the back of a horse | $(p = 1.3e - 03)$ |
| 3. A man riding on the back of a white horse | $(p = 0.9e - 03)$ |

The descriptions shown above have been extracted from the image of a profile belonging to the *scammer* category. The image shows a young man riding a brown horse. All three descriptions complement each other by adding additional details of the image. It is common to find misappropriated images that do not belong to the scammer—either because they have been stolen from a legitimate profile or because they have been taken from elsewhere on the Internet). A reverse search of the image does not reveal the source.

| Descriptions automatically generated from Fig. 4d: | |
|--|--------------------|
| 1. A man sitting in front of a laptop computer | $(p = 13.1e - 03)$ |
| 2. A man sitting at a table with a laptop | $(p = 3.6e - 03)$ |
| 3. A man sitting at a table with a laptop computer | $(p = 2.0e - 03)$ |

These descriptions belong to an image from the *scammer* category. The image shows a middle aged man sitting in front of a laptop and a table in the background. This image together with others found in the same profile are stock images.

It is worth noting the level of detail shown in each caption, which not only identifies the main actor within the picture (a man in these cases), but also the backdrop and the activity being undertaken.

There are a number of common topics displayed across images in both profiles. When looking at the gender of the people present in the images, we can observe that males appear in about 60% them as shown in Table 2. This matches with the distribution of gender reported in the profile demographics.

Table 2: Topics found across profiles with images.

| Type | Real Profiles | Scam Profiles | All |
|--------------|---------------|---------------|--------|
| Male | 57.75% | 63.76% | 60.48% |
| Groups | 0.50% | 2.22% | 1.28% |
| Children | 5.21% | 3.38% | 4.38% |
| Food | 1.86% | 3.62% | 2.66% |
| Animals | 0.77% | 1.08 % | 0.91% |
| Discriminant | 13.76% | 17.77% | 15.58% |

There are also a number of topics slightly more prevalent in one or the other profile categories, i.e.: group pictures (including couples), pictures with children, or presence of food (e.g., wine, bbq, cake). For instance, there are over four times more *group pictures* in the scammers category than in the real one. Contrastingly, the number of images with *children* in real profiles is almost double. Combining together all informative elements of the images, we can observe that about 15% of images contain descriptions that appear exclusively in one of the two categories (referred to as ‘*discriminant*’ profiles in Table 2). This indicates that there is a large number of images for which their context can be used to characterize scammers. In other words, scammer profiles feature more pictures of certain groups. Note that fraudsters frequently iterate through certain themes known to be appealing (e.g.: men in uniform). This might also simply be down to the availability of images which the scammer steals e.g., those from stock photo databases). The image shown in Fig. 4d, for instance, is a stock image.

2.4 Profile Descriptions

Contrary to most real-life encounters, on dating websites, a user can easily disclose very personal information, such as their life story, what they are looking for in a partner, their hobbies, their favorite music, etc., to a complete stranger and without being interrupted. Moreover, filling in a personal description is usually highly encouraged by any dating website, because it can capture other users’ attention and increase the chances of meeting a user’s ‘perfect match’.

For scammers, however, the profile description provides yet another means to mislead their victims. Prior research has shown that they will go to great lengths to create the ‘ideal’ profile, to gain a potential victim’s interest and to maintain the pretense of a real (online) relationship [31, 33]. As a result, most scammers in our dataset—5,027 out of 5,402—attempted to create an attractive user account by advertising broad pretended interests and characteristics. The real users were less inclined to provide

such personal information about themselves: only 5,274 (out of 14,720) generated a profile description.

Recent advances in natural language processing technology have enabled researchers to perform automatic linguistic analyses of lexical, morphological, and syntactic properties of texts. However, most traditional studies use large sizes of training data with a limited set of authors/users and topics, which usually leads to a better performance of the machine learning algorithms. Profile descriptions are, however, typically short and can include a whole range of different topics. With regard to the dataset described in this paper, the average number of words per profile description was 78.7, with scammers producing more words on average (104.5) than genuine users (54.1). This effect is so pronounced that despite there being fewer scammers than real users, the overall total of 525,336 words for the *scam* category was greater than the 285,407 words for the *real* category. The finding that scammers’ profiles have a higher word count compared to genuine profiles is consistent with previous literature stating that liars tend to produce more words [6]. To analyze the variety of topics that are present in our dataset, we used dictionary terms that are mapped to categories from the LIWC 2015 dictionary [18]. Category frequencies were recorded for each profile description. Our results showed that scammers referred considerably more to emotions—both positive and negative—than genuine users. Additionally, they use words related to family, friendship, certainty, males and females more often, while real users tend to focus on their motives or drivers (e.g., affiliation, achievement, status, goals), work, leisure, money, time and space. With regard to language use, we found that scammers use more formal language forms, while genuine users displayed more informal language forms (e.g., Netspeak).

3 Classifying False Profiles

A high-level overview of our system can be obtained from Fig. 5. The system is first trained using a dataset of real and scam profiles. The goal of this phase is to obtain the following key elements that will later be used to identify fraudsters:

- (i) A set of prediction models $\mathcal{P} = \{P_1, \dots, P_i\}$ that output the probability

$$\theta_i(\phi_1, \dots, \phi_n) = P_i[X = \text{scam} \mid (\phi_1, \dots, \phi_n)]$$

of each profile X being *scam* given a feature vector (ϕ_1, \dots, ϕ_n) obtained from different profile sections i .

- (ii) A weighted model $f(\mathcal{P}) = \sum w_i \cdot P_i$ that combines all individual predictions in \mathcal{P} . Here, each individual classifier P_i is weighted by w_i according to the accuracy given on a validation set that is different from the training one. This will also serve as a way to calibrate individual probabilities. The final classifier will then output a decision based on a vote such that

$$f = \begin{cases} \text{scam} & \text{if } f(\mathcal{P}) < \tau \\ \text{real} & \text{otherwise,} \end{cases}$$

where τ is a threshold typically set to $\left\lfloor \frac{\sum w_i}{2} \right\rfloor + 1$.

Table 3: Our proposed set of features.

| ID | Source | Name | Type | $ \theta_i $ |
|------------|--------------|--------------------------------|------|--------------|
| θ_M | Demographics | Age | NF | 237 |
| | | Gender | CF | |
| | | Latitude | NF | |
| | | Longitude | NF | |
| | | Country | CF | |
| | | Ethnicity | CF | |
| | | Occupation | CF | |
| θ_C | Captions | Marital Status | CF | 363 |
| | | $\text{set}(\text{entities})$ | CF | |
| | | $\text{set}(\text{actions})$ | CF | |
| θ_S | Descriptions | $\text{set}(\text{modifiers})$ | CF | 105,893 |
| | | $\text{set}(\text{ngrams})$ | SBF | |

For the sake of simplicity, we refer to the model presented in (ii) as *weighted-vote*. One can simplify the model by giving equal weight to all w_i (typically $w_i = 1$) and obtaining a nominal value for P_i before voting. In other words, applying a threshold for each P_i (e.g., 0.5) and creating an equal vote among participants. We refer to this non-weighted voting system as *simple-vote*.

3.1 Feature Engineering

Our work considers a diverse set of features in order to build a robust classification system. The proposed set of features contains elements obtained from three different sources: (i) structured attributes of the profile referred to as *demographics* and denoted as θ_M , (ii) features extracted from raw images referred to as *captions* (θ_C), and (iii) features extracted from unstructured text (*description* denoted as θ_S). Based on the preprocessing described in Section 2, we extract different types of features as described below.

- **Numerical Features (NF):** refers to those attributes from a profile that take a quantitative measurement such as the age (18-85) of a person.
- **Categorical Features (CF):** refers to those attributes that take a limited number of possible values such as the gender (male or female).
- **Set-based Features (SBF):** refers to those attributes that can take an arbitrary number of values and the relationship between sets of attributes is relevant (e.g., words in the description).

Table 3 shows the set of features proposed categorized by the source in the profile. For θ_M we considered the age, and the Cartesian location values as numerical values, while other attributes were treated as categorical. As described previously, common occupation responses were grouped into 45 different occupation areas (e.g., self-employed, military, legal). Long-tail occupation responses outside of these categories were grouped under *other*. In training, no such response appeared more than twice, and 85% of real and 96% of scammer occupations were captured by known categories.

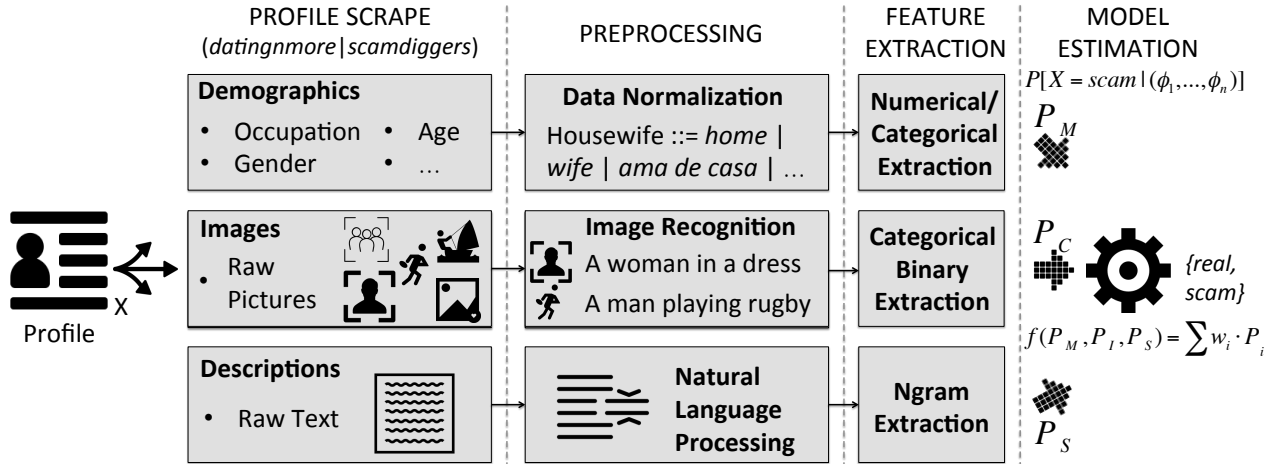


Figure 5: Key features extracted from dating profiles.

For θ_C , after extracting the most representative caption per image, we removed the least informative elements and retained only entities (nouns), actions (verbs), and modifiers (adverbs and adjectives). Each element in the caption was stemmed to a common base to reduce inflectional forms and derived forms. Rather than treating the captions as set-based features, we encoded them as categorical features, as the order of the actions is not relevant in this document vector approach. Generated captions are simple and their structure always follows the same pattern. The presence of a given action and the set of objects appearing in the image is itself informative. Encoding the relationship between the different parts of speech would unnecessarily increase the number of features.

Finally, we extracted set-based features from the textual content of the tokenized descriptions (θ_S). We considered Bag-of-Words (BOW), word n -grams, character n -grams and LIWC features. Because word bigrams (n -grams of length $n = 2$) yielded the best results during our preliminary experiments and combining different feature types did not lead to better classification results, we only included word bigrams in the rest of the experiments. Additionally, stemming and stop word removal resulted in a worse performance, so all features were included in their original form during the experiments.

3.2 Prediction Models

Because most of the fields are optional, user profiles in online dating sites are inherently incomplete. Some users are uncomfortable with high levels of self-disclosure and some are more interested in contacting others than presenting details about themselves [30]. Thus, any reliable detection system should be able to flexibly deal with incomplete profiles. In this section, we present three independent classifiers to estimate the presence of fraudulent profiles. Each classifier is designed to effectively model a section of the profile (based on θ_M , θ_C and θ_S as described previously). Probability outputs from each classifier are later combined to provide one balanced judgement. By using multiple classifiers designed on individual sections of the profile, we increase the likelihood that at least one classifier is capable of making an informed decision. Moreover, ensembles often perform better than single classifiers [3].

Demographics Classifier The demographics classifier uses the greatest variety of original profile attributes. Unlike the image and description classifiers, handling this feature set means dealing with non-binary missing data situations—location and ethnicity might be missing for a given profile which still contains age and gender information. When no data is available, the least-informative prior is the base rate of real vs scam profiles, as is used in the image and description classifiers. In most situations within the demographics data more information than this is available, and should be used.

Approaches for handling missing data—in the situation where the case cannot be discarded, such as appearing in a test or validation set—reduce to either some form of imputation or the use of a classifier robust to missing values, such as a Naive Bayes approach. Given problematic randomness assumptions inherent to the most useful imputation methods, we opt to use a Naive Bayesian classifier to handle prediction for profiles with missing data attributes.

However, Naive Bayes is not the most effective classifier for profiles with all data present—a significant proportion of the dataset. For this subset, it is more effective to use a classifier which performs better and does not handle cases with missing data. In our case, a Random Forests model was selected. The final approach to providing $P_M(X = scam)$ is to train a joint Random Forests and Naive Bayes model, using the high-performing Random Forests model to make predictions where all demographic data is available, and the gracefully-degrading Naive Bayes model for all other cases.

Images Classifier We build a prediction model based on the features extracted from the captions of the images such that $P_C(X = scam)$. The architecture of our system is highly flexible and accepts a wide range of classifiers. Our current implementation supports Support Vector Machine (SVM), Random Forests and Extra Randomized Trees. For the purpose of this paper, we selected SVM with radial kernel as the base classifier for the images.

SVM has been successfully applied to fraud detection in the past [1] and has been shown to have better performance compared to 179 classifiers (out of 180) on various datasets [5]. SVM

also tends to perform well when the number of samples is much greater than the number of features, as it is the case here. In addition, it is less sensitive to data outliers—instead of minimizing the local error, SVM tries to reduce the upper bound on the generalization error.

Descriptions Classifier Previous work [33] has shown that scammers attempt to keep their labor costs down to be able to exploit different social media and to continuously produce the interaction that is required to make their on-going scams succeed. To achieve this goal, they tend to edit pre-written scripts that are often shared on underground forums—labeled by the ethnicity, age group, location and gender of the potential victim. Hence, for providing $P_S(X = \text{scam})$, we compared the performance of two approaches: (i) a similarity-based approach, in which we applied shingling ($k = 5$) to extract the set of all substrings from each profile description in training and calculated the Jaccard similarity for each pair of profile descriptions (see [14]); and (ii), we trained an SVM algorithm (linear kernel) as implemented in LibShortText [34], an open-source software package for short-text classification and analysis. Parameters for both approaches were experimentally determined on a small subset of each training partition during cross validation. Within the SVM experiments, features were represented by TF-IDF scores, which reflect the importance of each feature (in this case, each word bigram) to a document (i.e. a user profile description) in terms of a numerical frequency statistic over the corpus [14].

Ensemble Classifier The goal of this method is to combine the predictions of the base estimators described above to improve the robustness of the classification. Ensemble methods are designed to construct a decision based on a set of classifiers by taking a weighted vote of all available predictions. In our system, we have a function f that is estimated using an independent set of samples. This function will then be used during testing to weight each prediction model P_i such that: $f(P_M, P_C, P_S) = \{\text{scam}, \text{real}\}$.

For the decision function f we use a Radial Basis Function (RBF) that measures the distance to the center of the SVM hyperplane bounding each P_i . This function is defined on a Euclidean space and it only measures the norm between that point and the center (without considering the angular momentum). This function is approximated with the following form

$$f = \sum w_i \delta(\|p_i\|),$$

which can be interpreted as the sum of the weights w_i times the probability score $p_i \in P_i$ given by the individual classifiers in the voting system described above.

Single Classifier We compare the results of our ensemble method to the predictions made by a single SVM classifier (linear kernel) in which all demographics, captions and description features are included in each document instance. Features were represented by their absolute values and parameters were again experimentally determined on a small subset of each training partition during cross validation.

4 Evaluation

In evaluating and developing the classification system described before, we applied the following methodology.

Methodology We divided the dataset into a 60% training set, a 20% test set and a 20% validation set. Profiles were assigned to each set randomly under a constraint preventing variants of the same scam profile from being assigned different sets or folds. Development of the classification system proceeded as follows:

1. Each component classifier was designed within the 60% training set, and individual performance levels established through ten-fold cross-validation within this set.
2. Once classifier design was complete, each component classifier was trained on the full training set.
3. To design the ensemble model, each classifier produced probabilities and labels for the test set. The ensemble was developed on these probabilities, and performance was established through five-fold cross-validation on the test set.
4. Based on performance within the test set, the ensemble model and the choice of outcomes to report within the validation set was decided.
5. For final validation, individual classifiers were trained on the training set, produced probabilities and labels for the testing set and the validation set, the ensemble model was trained on the probabilities given for the test set, and its predictions taken for the validation set.
6. The single classifier was trained on the combination of the training and test data and evaluated on the validation set.

4.1 Classification Results

We present our results together with a number of case studies, covering all four dimensions of the classification performance: (i) scam profiles correctly classified (TP), (ii) real profiles correctly classified (TN), (iii) real profiles misclassified (FP), and (iv) scam profiles misclassified (FN).

Summary Table 4 presents the results within the validation set for each classifier, for simple majority voting between all three classifier outputs, and for the SVM ensemble model trained on the classifier probabilities given for the test set. Precision, recall and F1 are given for predicting scam profiles (the minority class). Judging performance by F1, the best individual classifier was the SVM description classifier (F1 = 0.842). As can be expected, the similarity-based approach (threshold Jaccard similarity of 0.259) yielded a high precision score, but a low recall score, which resulted in a markedly lower F1 of 0.712. The demographics classifier was the next best component classifier (F1 = 0.840), but the captions classifier was highly precise, making only two false-positive judgements. Simple majority voting between classifier labels improved performance significantly compared to any individual classifier, raising F1 to 0.904, with a precision of 0.996. A single classifier using all features outperformed majority voting (F1 = 0.927). The ensemble system outperformed both the

Table 4: Final results for each component classifier, simple majority voting, a similarity-only approach, a single classifier using all features, and the weighted-vote ensemble

| CLASSIFIER | TN | FN | FP | TP | PREC. | REC. | F1 | ACC |
|------------------------|------|-----|-----|------|-------|-------|-------|-------|
| demographics | 2725 | 196 | 149 | 903 | 0.858 | 0.822 | 0.840 | 0.913 |
| captions | 2872 | 499 | 2 | 600 | 0.997 | 0.546 | 0.705 | 0.874 |
| description | 2758 | 215 | 116 | 884 | 0.884 | 0.804 | 0.842 | 0.917 |
| similarity-only | 2939 | 435 | 28 | 571 | 0.953 | 0.568 | 0.712 | 0.884 |
| simple-vote | 2870 | 189 | 4 | 910 | 0.996 | 0.828 | 0.904 | 0.951 |
| single | 2820 | 108 | 54 | 1027 | 0.950 | 0.905 | 0.927 | 0.959 |
| weighted-vote | 2834 | 78 | 40 | 1021 | 0.962 | 0.929 | 0.945 | 0.970 |
| Excluding new variants | | | | | | | | |
| demographics | 2725 | 122 | 149 | 569 | 0.792 | 0.823 | 0.808 | 0.924 |
| captions | 2872 | 378 | 2 | 313 | 0.994 | 0.453 | 0.622 | 0.893 |
| description | 2758 | 119 | 116 | 572 | 0.831 | 0.828 | 0.830 | 0.934 |
| simple-vote | 2870 | 129 | 4 | 562 | 0.993 | 0.813 | 0.894 | 0.963 |
| weighted-vote | 2818 | 53 | 56 | 638 | 0.919 | 0.923 | 0.921 | 0.969 |
| Excluding all variants | | | | | | | | |
| demographics | 2707 | 114 | 167 | 577 | 0.776 | 0.835 | 0.804 | 0.921 |
| captions | 2874 | 426 | 0 | 265 | 1.000 | 0.384 | 0.554 | 0.881 |
| description | 2731 | 171 | 143 | 520 | 0.784 | 0.753 | 0.768 | 0.912 |
| simple-vote | 2860 | 159 | 14 | 532 | 0.974 | 0.770 | 0.860 | 0.951 |
| single | 2829 | 98 | 45 | 592 | 0.929 | 0.858 | 0.892 | 0.960 |
| weighted-vote | 2841 | 69 | 33 | 622 | 0.950 | 0.900 | 0.924 | 0.971 |

single classifier and majority voting at 0.945 F1, significantly improving recall whilst maintaining a high level of precision. Over 97% of all profiles were classified correctly. Fig. 6 characterises the ROC performance for this ensemble depending on whether variants (near-duplicate profiles) were excluded.

Feature Analysis We describe some of the most important features as identified by our classifiers.

Table 5a presents the total decrease in node impurities from splitting on the each feature in the RF component of the demographics model, averaged over all trees. The most important feature was the occupation area reported in the profile. Node purity rankings are known to bias towards factors with many levels, but the size of the interval between the occupation area and the other features suggests that this ranking is genuine. This would accord with our observations in 2.2 about the use of occupation area as an attractive status marker.

Table 5b presents the highest-weighted bigrams from the descriptions classifier for the purpose of predicting the scam category. The most informative features tend to relate to nonfluencies in English (starting descriptions with ‘Im’ or ‘Am’, constructs like ‘by name’) and attempts to overtly signal a romantic or trustworthy nature (e.g., “caring”, “passionate”, “loving”). Our topic analysis in Section 2.4 also captures this tendency of scam profiles to include more emotive language.

Table 5c presents the most discriminant features for the captions classifier. Features with a negative weight are more informative when discriminating real profiles. Instead, features with positive weight relate to scam profiles. Interestingly, some of the top elements embedded in the images map with relevant traits observed in the demographics such as the occupation (e.g., military) or the gender (e.g., male) c.f. §2.2.

True Positives About 98% of the scam profiles have been detected by at least one of the classifiers. Consensus between classifiers accounts for the majority of TPs, but performance

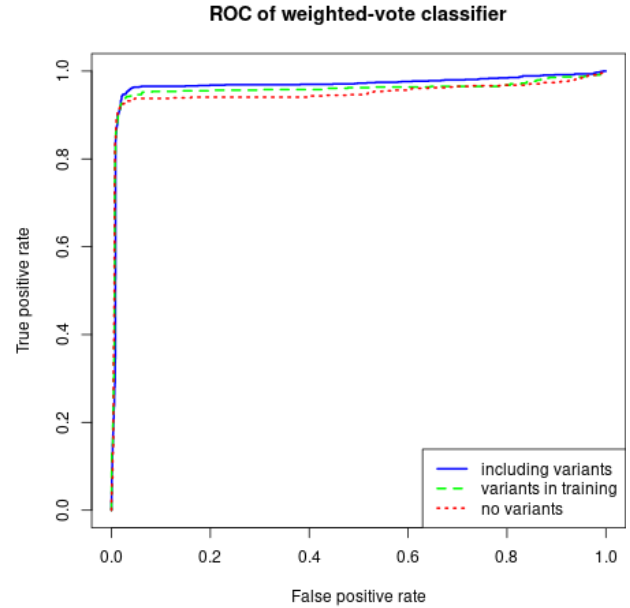


Figure 6: ROC for the ensemble classifier

Table 5: Top-ranked features for component classifiers

| (a) Feature | ranking | de- | (b) Top-weighted bigrams for | (c) Feature raking captions | |
|-------------|---------|------------|------------------------------|-----------------------------|--------|
| mographics | RF | scam | descriptions | captions | |
| feature | purity | bigram | weight | Keyword | weight |
| occupation | 332.70 | <start> im | 0.3086 | pizza | -1.0 |
| latitude | 198.02 | don t | 0.2318 | picture | -0.52 |
| status | 128.76 | caring and | 0.1776 | child | -0.50 |
| longitude | 128.71 | and caring | 0.1674 | bottle | -0.46 |
| age | 114.24 | by name | 0.1644 | christmas | -0.46 |
| ethnicity | 110.53 | <start> am | 0.1643 | driving | 1.0 |
| gender | 64.52 | am just | 0.1641 | military | 1.0 |
| | | that will | 0.1572 | birthday | 2.0 |
| | | am here | 0.1568 | group | 2.46 |
| | | tell you | 0.1481 | male | 2.95 |

improves yet further when resolving disputes using classifier weights learned on an independent sample in the ensemble voting scheme. Under this scheme, we manage to detect about 93% of the fraudulent profiles with a high degree of confidence, compared to 81% when only relying on the simple voting scheme. Roughly 36% of scammers were identified by all three classifiers.

For illustration, we present one TP case randomly chosen from those identified with a high degree of confidence. This is the case of a profile presenting as a 26 year old African American female, with the occupation reported as “student”. In the description, we can see certain traits uncommon in legitimate profiles, like the intention to establish a “great friendship”. The misprint on the occupation and the poor language proficiency in the description might indicate that the fraudster lacks fluency in the English language.

Hi., Am Vivia. How are you doing? I would be very happy to have a great friendship with you. my personal email is (<user>@hotmail.com) I look forward to hearing from you, Vivia

There are two images in the profile, which match the demographics reported. Specifically, one of the images has a young

woman sitting on a park bench, while the other shows the same woman sitting in a study room with a laptop. The prevalence of certain elements such as the use of laptops across stock scam profile images could explain the decision taken by the image classifier.

True Negatives All real profiles have been identified as such by at least one of the classifiers. When combining the decisions, our system correctly classified 99.9% real profiles using *simple-vote* and about 98.6% using *weighted-vote*.

Our randomly selected exemplar case is that of a 55 year old white woman. This user is based in Mexico, a location comparatively underpopulated by scam profiles. It is worth noting that the age also deviates significantly from the average age of female scammers (30).

This profile had an avatar image showing the face and shoulders of a woman. All the elements in the image looked conventional, which is most likely why the classifier identified the profile as real. It is worth noting the poor quality of the image. Although the quality of the images was not measured in this work, we noticed that fraudsters care about it, and intend to investigate this further. Comparing the profile description to the previous TP example, the overall fluency is notably higher, both in terms of English grammar and the appropriate format (as a self-description rather than a message). The user describes herself and her interests, rather than focusing entirely on the reader.

False Positives There are a total of 4 real profiles misclassified under our most precise setting, for a false positive rate of 0.1%. For those 4 profiles there was at least one classifier that correctly predicted the profile and, interestingly, all three classifications were available. This means that none of the errors can be attributed to a lack of information. Under weighted-voting the false positive rate rises to 1.4%. When looking at the errors common to both voting schemes, we only find 2 misclassified real profiles. Both errors come from predictions by the demographics and descriptions classifiers.

From the demographics, we see that one profile is widowed, common amongst scammers, and the other is mixed-race, which is slightly more common amongst scammers. Both users claim locations with high scammer population: Texas in the US, and close to London in the UK. One of the profile descriptions was very short, and focused on the nature of their intended relationship and the qualities of the intended partner. The other was long, but both referred to topics which are more strongly correlated with fraud, such as relationships and positive emotions.

False Negatives For misclassified scam profiles, weighted-voting produced 78 errors and simple-voting produced 189 errors. Out of all errors, we find 22 cases where all three individual classifiers failed. Manual analysis of these cases reveals that certain parts of the profiles look genuinely normal. In general, the image caption classifier was most likely to produce false negatives, and overall errors occurred when either of the demographics or description classifiers agreed. Of the two, the description classifier was slightly more likely to produce FNs.

Such is the case of a 45 year old American soldier from Louisiana named Larry. While the description is ordinary, the

occupation raises suspicions due to the prevalence of military scammers. When there is a technical draw between the demographics and the descriptions and the images are not informative enough one can only aim at extracting additional features. For instance, in the case of Larry, the messages he exchanged with a victim are valuable. The excerpt below shows a message sent to one of his victims wherein one may observe distinguishable wording and a clear manipulative strategy in parts highlighted below:

I must confess to you, you look charming and from all I read on your profile *Id want you to be my one special woman*. I wish to build a one *big happy family* around you. Im widowed with two girls, Emily and Mary, I lost their mom some years ago and since then *Ive been celibate*. *I think youve got all it takes to fill the vacuum left by my late wife to me and the kids*. I seek to grow old with you, children everywhere and grey hairs on our head. *I wanna love you for a life time*. Hope to read from you soon. You can addme on facebook <user> or mail me at <user>@yahoo.com. Looking forward to your message. Regards, Larry. —A message sent to a victim

Adding such elements from later steps in the scam life-cycle could help to cover the cases where all three individual classifiers failed. However, these features are not necessarily available for all deployment situations, and would restrict more general application of the model.

5 Discussion

In this paper, we established the need for a systematic approach that automates online romance scammer detection. The system we presented is a key first step in developing automated tools that are able to assist both dating site administrators and end-users in identifying fraudsters *before* they can cause any harm to their victims. However, there are risks in deploying such automated systems, principally:

- i) Risk of denying the service to legitimate users.
- ii) Risk of having scammers that have evaded the system.

We next discuss these implications and the limitations.

5.1 Context-based Performance Maximization

Under our current system, 96% of profiles identified as scammers truly are, and about 93% of all scammers are detected. This performance is optimized for the harmonic mean of these rates. One might further tune the model, either for minimizing the detection of *false positives* or minimizing the *false negatives*. This decision will rely on the priorities of the user of a classification tool:

Minimizing FP – when real profiles are misclassified, users are inconvenienced by being flagged as scammers and are likely to be annoyed at a platform that does this. Thus, detection systems being run by dating sites must review alerts or risk losing customers. To reduce workload and costs, dating sites may want to minimize the risk of misclassifying real users, and use user-reporting and education tools to catch scammers who evade preemptive detection.

Minimizing FN – when scam profiles are misclassified, a user risks being exposed to a scammer and suffering emotionally

and financially as a result. Given that the opportunity cost is comparatively low for potential partners being filtered out, a “better safe than sorry” attitude is justified. As such, safe-browsing tools that a user deploys themselves may wish to bias towards always flagging scammer profiles, as the user may always disable or ignore such a tool if convinced it is in error⁵.

The simple voting classifier presented in Table 4 provides an easy example of a system biased towards minimising false positives, with only 4 appearing in a set of nearly 3,000 real profiles. If all three classifiers were required to agree before a profile was classified as a scam (unanimous voting), then the false positive rate would be too low for this study to detect (0 observed). Alternately, if the firing of any classifier was sufficient reason to flag a profile, only 22 scammer profiles in over 1,000 would escape being flagged.

Returning to the better-performing machine-weighted voting system, the ensemble system could be optimized for any risk ratio by optimization towards a modified F-score. The F-score can be weighted towards any desired ratio of *precision* (minimizing false positives) and *recall* (minimizing false negatives) by adjusting the β parameter in the general equation:

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

where the β expresses the ratio by which to value recall higher than precision. By selecting an appropriate balance between these measures and then evaluating classifiers against this measure, a weighted voting system can be tuned to individual risk tolerances.

5.2 Comparison with Moderator Justifications

The moderators who identified profiles as romance scammers provide a list of justifications for their decision on each profile. By analyzing the given justifications, we can examine our classifier’s performance next to individual human strategies for scammer identification.

Table 6 presents figures for the proportion of scam profiles labeled with common justifications. The figures are counted in terms of profiles and not profile-variants. Alongside figures for all scam profiles are figures for the validation set upon which the ensemble system was tested, and figures for the scam profiles which the ensemble classifier mislabeled as non-scam profiles (false negatives).

Certain observations can be made. Firstly, on overall justification proportions across scam profiles, we can see that examination of the geolocation of a scammer’s IP address is a heavily relied-on method for moderators, with contradictions between this and the profile’s stated location being a justification listed for 87% of all scam profiles. The next most common justification was that a profile uses suspicious language in its self-description: expressions of this ranged from identification of “Nigerian wording” to moderators recognizing text being reused from previous scams .

⁵Such tools may of course also allow the user to define their own risk tolerances.

Table 6: Comparison of overall, validation and false-negative incidence of moderator justifications for scam-classified profiles

| REASON | ALL SCAMS | VALID. | FN | REC. |
|-----------------------------|------------|-----------|----------|------|
| IP contradicts location | 3030 (87%) | 620 (87%) | 44 (85%) | 0.93 |
| Suspicious language use | 2499 (72%) | 507 (71%) | 34 (65%) | 0.93 |
| IP address is a proxy | 2156 (62%) | 433 (60%) | 25 (48%) | 0.94 |
| Known scammer picture | 1379 (40%) | 299 (42%) | 17 (33%) | 0.94 |
| Known scammer details | 1368 (39%) | 284 (40%) | 13 (25%) | 0.95 |
| Self-contradictory profile | 1145 (33%) | 242 (34%) | 12 (23%) | 0.95 |
| IP location is suspicious | 968 (28%) | 211 (29%) | 22 (42%) | 0.90 |
| Mass-mailing other users | 761 (22%) | 168 (23%) | 10 (19%) | 0.94 |
| Picture contradicts profile | 261 (7%) | 55 (8%) | 4 (8%) | 0.93 |

Comparison of proportions between the overall dataset and validation set show little deviation in justification proportion, demonstrating a lack of bias. By comparing proportions within the false-negative profiles to those in the overall validation set, we may discern any systemic differences in identification rate.

Most justifications show similar or lower proportions in the false-negative profiles, indicating that the ensemble is either no worse than average within these subcategories, or may be better than average. One category of justifications alone showed worse performance for the ensemble—where the human moderators judged that the IP-determined origin of the scammer was in a country they deem suspicious (e.g., a West African nation). The recall of profiles justified with this reason was 0.9, lower than average. IP address information is not available for non-scam users in our dataset, so this discrepancy cannot be fully investigated, but it might suggest that the partially location-based demographics classifier is not yet matching expert understanding of scam-correlated locations.

5.3 Evasion

From the previous section, we see that moderators heavily rely on certain features, such as the IP address, that can easily be obscured. Moderators also check if the IP address is known to belong to a proxy. Although this could certainly be used as a feature in our system, scammers could use unregistered proxies (as yet inexpensive) or compromised hosts in residential IP address space to evade detection.

In contrast, our system relies on a wide range of features that are more difficult to evade, such as textual features⁶; or features for which their obfuscation might render a profile unattractive, such as the demographics or the images.

A natural scammer response to this kind of profile detection could be to cease hand-crafting unusually attractive or targeted profiles, and instead turn to cloning the existing profiles of real users from different dating sites. By preferentially cloning attractive profiles, they would retain the high match-rate which enables contact with potential victims, and by using real users’ profile information they would avoid detection systems geared towards their own idiosyncratic profile elements. Scammers are already partially engaging in this sort of behaviour when they re-use images of real people taken from the web.

⁶Prior work in computational linguistics has shown that the combination of (un)consciously made linguistic decisions is unique for each individual — like a fingerprint or a DNA profile [27, 20]

Table 7: Comparison of the profile elements used in our classification experiments with availability of these elements on popular dating sites. (✓: present; %: requires inference)

| SITE | age | gender | ethn. | marital | occ. | location | image | descr. |
|-----------------------|-----|--------|-------|---------|------|----------|-------|--------|
| <i>datingmore.com</i> | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| <i>match.com</i> | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| <i>okcupid.com</i> | ✓ | ✓ | ✓ | ✓ | % | ✓ | ✓ | ✓ |
| <i>pof.com</i> | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| <i>eharmony.com</i> | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| <i>tinder.com</i> | ✓ | ✓ | | | | % | ✓ | ✓ |

The solution to such a development will rely on the deployment of profile-cloning detection systems, such as that described by Kontaxis et al. [12], perhaps augmented by behavioural classifiers operating on e.g., the language used in messages.

5.4 Comparison of Dating Site Elements

The *datingmore.com* site used as the source of data for our experiments is comparatively small, and has a niche appeal due to its intensive moderation by experts in identifying online dating fraud. It is therefore worthwhile considering its comparability to other dating sites, as a first step to understand the generalisability of our results.

Table 7 compares the features from the *datingmore.com* profiles which were used in our classifier with the profile elements available on five market-leading dating sites. Coverage was good. All three of the ensemble components would be able to operate across these sites. The image and description profile elements are always supported, and at least some demographic information is always available. The dating site with the fewest profile elements in common with our features is *tinder.com*, which has a distinctive locality-based use case which may hinder the online dating fraud our system aims to detect.

5.5 Limitations & Deployment Considerations

There are limitations to our work which must be borne in mind. Firstly, whilst we have taken pains to make use of profile features which should be visible on other dating platforms, we have not yet tested our classification approach on profiles from other dating sites. It may be the case that scammers and/or real users show different characteristics in different dating platforms, which would limit the applicability of our method. We are currently seeking other sources of user profile data to investigate this possibility. More information on scammer/user traits could generally inform ongoing research into—and prevention of—online dating fraud.

Secondly, our results show a number of false negative classifications. Further inspection of the data on scammers suggests that augmenting our approach with other classifiers—such as ones using geolocated IP addresses or observations of on-platform behavior (e.g., messaging)—could help capture these scammers where the public profile information is inconclusive.

Online dating sites could deploy a detection system such as ours on their premises at the profile registration stage. Security administrators would then be responsible for validating and acting upon the output of the classification system. How dating sites can responsibly anticipate and respond to errors in automated classification systems such as this is a point of policy on which

practical industry insight would be highly valuable. More generally, this raises the issue of accountability under the EU General Data Protection Regulation, and the “right to explanation” of algorithmic decisions that significantly affect a user.

In the case where our system is deployed locally, the implications of our decisions can have a paramount effect on the user. Suppose that our system predicts that a given scam profile is safe. This might give the user a false sense of security, and encourage them into beginning a relationship with less caution than they would have applied otherwise. Designing a tool which protects users while minimizing the risk of blind trust is a challenging interface design problem, but one which is outside the scope of this paper.

6 Related Work

Despite the rapidly increasing number of victims⁷, previous work on online dating fraud is limited, focussing mainly on case studies [21], interviews with online dating site users about their security practices [17] and interview and questionnaire-based psychological profiling of victims [2, 31, 32].

Three recent studies provide insight into romance scammer strategies. The authors of [33] carried out a study on the personals section of Craigslist to identify common methods of romance scammers responding to honeypot adverts. Alongside identifying approaches outside of traditional trust-building relationship-oriented scammers (including driving users to other platforms and delivering hooks for premium-rate numbers), they observed scammers were mostly West African in origin, and used scraped images of attractive women [33]. Huang et al. [9] performed a large-scale study of dating profiles that are used by scammers, covering 500,000 scam accounts from an anonymous Chinese online dating site. They found that different types of scammers target different audiences and that advanced scammers are more successful in attracting potential victims’ attention. Finally, the authors of [4] describe geographic variation in dating fraud profiles, and propose a set of methods to improve geolocation when attackers are hiding behind proxies. While assisting in attribution of origins, these methods cannot be used for scammer detection.

Although a number of solutions based on machine learning techniques already exist to detect malicious activity on online services (e.g., detection of spam [13, 25, 29], or false identities [20, 15]), to our knowledge, no prior work has attempted to automatically detect romance scammers. One of the main reasons for this is that the dynamics of dating websites make scam detection more difficult than in other domains, such as email or social networking. The intended operation of a dating site is that previously unconnected users will reach out and initiate contact with people they do not know, and so spontaneous, unsolicited communications cannot be viewed as a reliable signal of malicious behavior. Activities that in other areas might be considered suspicious—contacting many users, providing false profile attributes, migrating conversations to other media—could all also be considered normal behavior amongst dating site users [9]. Moreover, romance scams are for the most part carried out by

⁷According to research at the Chartered Trading Standard Institute in the UK, the number of romance scam victims will more than triple by 2019 [16].

humans, adapting to changing circumstances, and so approaches which rely on detecting bot-like behavior are similarly stymied. In this paper, we address these issues by analyzing the launching point for the scam—the user profile.

7 Conclusions

In this paper, we have presented the first framework systematizing the identification of false online dating personas. Our exploratory analysis identified the sugarcoated lures used on fraudulent profiles. By analyzing the prevalence of these traits with respect to legitimate profiles, we engineered a diverse discriminatory feature set, using state-of-the-art text and image processing from multiple profile segments. This feature set allowed us to develop a set of independent classification systems which adjust to the omission of profile details.

Our experimental results show that our system can accurately detect online dating fraud profiles, with high precision. A case by case analysis of our results, however, indicates that there are certain false profiles that look genuinely real. For these cases, we have noted that other sources of information, such as the messages exchanged, could be very informative. As future directions, we aim to more broadly examine the available data on online dating fraud, seeking information actionable for enforcement and other countermeasures. We also hope to explore the question of how, at a local level, interventions designed to warn and protect users from scammers can avoid forming dependences that reduce awareness.

Acknowledgements

This work is supported by award EP/N028112/1 “DAPM: Detecting and Preventing Mass-Marketing Fraud (MMF)”, from the UK Engineering and Physical Sciences Research Council. This research would not be possible without the work of the operators and moderators of the scamdigger.com and datingnmore.com websites, and their commitment to transparently combatting online dating fraud.

Availability

To enable replication and foster research in this area we release the code used to obtain our data, the processing steps taken to prepare it, and the implementations of each classifier together with additional details of our results, all available online at <https://github.com/gsuareztangil/automatic-romancescam-digger>.

References

- [1] S. Bhattacharyya, S. Jha, K. Tharakunnel, and J. C. Westland. Data mining for credit card fraud: A comparative study. *Decision Support Systems*, 50(3):602–613, 2011.
- [2] T. Buchanan and M. T. Whitty. The online dating romance scam: causes and consequences of victimhood. *Psychology, Crime & Law*, 20(3):261–283, 2014.
- [3] T. G. Dietterich. Ensemble methods in machine learning. In *International Workshop on Multiple Classifier Systems*, pages 1–15. Springer, 2000.
- [4] M. Edwards, G. Suarez-Tangil, C. Peersman, G. Stringhini, A. Rashid, and M. Whitty. The geography of online dating fraud. In *Workshop on Technology and Consumer Protection (ConPro)*, 2018.
- [5] M. Fernández-Delgado, E. Cernadas, S. Barro, and D. Amorim. Do we need hundreds of classifiers to solve real world classification problems? *The Journal of Machine Learning Research (JMLR)*, 15(1):3133–3181, Jan. 2014.
- [6] J. T. Hancock, L. Curry, S. Goorha, and M. Woodworth. Automated linguistic analysis of deceptive and truthful synchronous computer-mediated communication. In *System Sciences, 2005. HICSS’05. Proceedings of the 38th Annual Hawaii International Conference on*, page 22c. IEEE, 2005.
- [7] J. T. Hancock, C. Toma, and N. Ellison. The truth about lying in online dating profiles. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 449–452. ACM, 2007.
- [8] D. Hembree. Online Romance Scams Are Fleecing More Americans. <https://www.forbes.com/sites/dianahembree/2017/06/20/romance-scam-crimes-on-the-rise/>, 2017. Online; accessed October 2017.
- [9] J. Huang, G. Stringhini, and P. Yong. Quit playing games with my heart: Understanding online dating scams. In *Proceedings of the International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment*, pages 216–236. Springer, 2015.
- [10] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137, 2015.
- [11] D. T. Kenrick, E. K. Sadalla, G. Groth, and M. R. Trost. Evolution, traits, and the stages of human courtship: Qualifying the parental investment model. *Journal of Personality*, 58(1):97–116, 1990.
- [12] G. Kontaxis, I. Polakis, S. Ioannidis, and E. P. Markatos. Detecting social network profile cloning. In *Pervasive Computing and Communications Workshops (PERCOM Workshops)*, 2011 *IEEE International Conference on*, pages 295–300. IEEE, 2011.
- [13] K. Lee, J. Caverlee, and S. Webb. Uncovering social spammers: social honeypots+ machine learning. In *Proceedings of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 435–442. ACM, 2010.
- [14] J. Leskovec, A. Rajaraman, and J. D. Ullman. *Mining of massive datasets*. Cambridge University Press, 2014.

- [15] W. Magdy, Y. Elkhatib, G. Tyson, S. Joglekar, and N. Sasstry. Fake it till you make it: Fishing for catfishes. *arXiv preprint arXiv:1705.06530*, 2017.
- [16] K. Morley. One million pensioners will be on ‘suckers lists’ by 2019, 2017.
- [17] B. Obada-Obieh, S. Chiasson, and A. Somayaji. ‘Don’t break my heart!’: user security strategies for online dating. In *Proceedings of the Usable Security Mini Conference (USEC)*. Internet Society, 2017.
- [18] J. W. Pennebaker, R. L. Boyd, K. Jordan, and K. Blackburn. The development and psychometric properties of LIWC2015. Technical report, 2015.
- [19] D. Ramalingam and V. Chinnaiyah. Fake profile detection techniques in large-scale online social networks: A comprehensive review. *Computers & Electrical Engineering*, 2017.
- [20] A. Rashid, A. Baron, P. Rayson, C. May-Chahal, P. Greenwood, and J. Walkerdine. Who am I? analyzing digital personas in cybercrime investigations. *Computer*, 46(4):54–61, 2013.
- [21] A. Rege. What’s love got to do with it? exploring online dating scams and identity fraud. *International Journal of Cyber Criminology*, 3(2):494, 2009.
- [22] H. Roberts. There are more ‘romance scam’ victims than ever – and people were defrauded out of 39 million in 2016. *businessinsider.com*: <https://goo.gl/P9Bj6M>, 2016. Online; accessed April 2018.
- [23] RomanceScam. Romance Scam: IP addresses of scammers. <https://www.romancescam.com/ipsearch.php>, 2013. Online; accessed Oct 2017.
- [24] ScamWatch. Scam statistics. <https://www.scamwatch.gov.au/about-scamwatch/scam-statistics>, 2017. Online; accessed October 2017.
- [25] G. Stringhini, C. Kruegel, and G. Vigna. Detecting spammers on social networks. In *Proceedings of the 26th Annual Computer Security Applications Conference*, pages 1–9. ACM, 2010.
- [26] L. Tung. FBI: Victims of online fraud lost \$800m to scammers last year, 2015.
- [27] H. Van Halteren, H. Baayen, F. Tweedie, M. Haverkort, and A. Neijt. New machine learning methods demonstrate the existence of a human stylome. *Journal of Quantitative Linguistics*, 12(1):65–77, 2005.
- [28] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: Lessons learned from the 2015 MSCOCO image captioning challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2016.
- [29] A. H. Wang. Don’t follow me: Spam detection in twitter. In *Security and Cryptography (SECRYPT), Proceedings of the 2010 International Conference on*, pages 1–10. IEEE, 2010.
- [30] M. T. Whitty. Revealing the ‘real’ me, searching for the ‘actual’ you: Presentations of self on an internet dating site. *Computers in Human Behavior*, 24(4):1707–1723, 2008.
- [31] M. T. Whitty. The scammers’ persuasive techniques model: Development of a stage model to explain the online dating romance scam. *British Journal of Criminology*, 53(4):665–684, 2013.
- [32] M. T. Whitty and T. Buchanan. The online dating romance scam: The psychological impact on victims—both financial and non-financial. *Criminology & Criminal Justice*, 16(2):176–194, 2016.
- [33] T.-F. Yen and M. Jakobsson. Case study: Romance scams. In *Understanding Social Engineering Based Scams*, pages 103–113. Springer, 2016.
- [34] H. Yu, C. Ho, Y. Juan, and C. Lin. Libshorttext: A library for short-text classification and analysis. *Technical report*, 2013.