

# Identifying tweets from Syria refugees using a Random Forest classifier

Smarti Reel  
Health Informatics Centre  
University of Dundee  
Dundee, UK  
s.reel@dundee.ac.uk

Patrick Wong  
School of Computing and  
Communications  
Open University  
Milton Keynes, UK  
patrick.wong@open.ac.uk

Belinda Wu  
School of Politics,  
Philosophy, Economics,  
Development, Geography  
Open University, UK  
belinda.wu@open.ac.uk

Soraya Kouadri Mostefaoui  
School of Computing and  
Communications  
Open University, UK  
soraya.kouadri@open.ac.uk

Haiming Liu  
Department of Computer  
Science and Technology  
University of Bedfordshire,  
Luton, UK  
Haiming.Liu@beds.ac.uk

**Abstract**— A social unrest and violent atmosphere can force a vast number of people to flee their country. While governments and international aid organizations need migration data to inform their decisions, the availability of this data is often delayed due to the tediousness to collect and publish this data. Recent studies recognized the increasing usage of social networking platforms amongst refugees to seek help and express their hardship during their journeys. This paper investigates the feasibility of accurately extracting and identifying tweets from Syria refugees. A robust framework has been developed to find, retrieve, clean and classify tweets from Syria. This includes the development of a Random Forest classifier, which automatically determines which tweets are from Syria refugees. Testing the classifier with samples of historical Twitter data produced promising result of 81% correct classification rate. This preliminary study demonstrates the potential that refugees' messages can be accurately identified and extracted from social media data mixed with many unwanted messages, and this enables further works for studying refugee issues and predicting their migration patterns.

**Keywords**— *Social media, Syria, refugees, classification*

**Short Papers, CSCSI-ISNA**

## I. INTRODUCTION

When a large-scale disaster or conflict breaks out, vast numbers of people in the affected areas often migrate and seek refuges in safe havens elsewhere. Such huge sudden movement of people can cause complex humanitarian, social and economic issues to the refugees, the nearby and hosting countries as well as the various governmental and aid agencies. If the patterns of these migrations can be correctly understood and accurately predicted, the governments and aid agencies of the countries concerned can be better prepared to help the refugees. Recent studies have recognized the increasing use of social media such as Twitter and Facebook amongst refugees/migrants [1]. They often use social networking platforms to express the hardship and difficulties they face and report their experience, as well as using such platforms as a tool in migration planning and finding out information during their actual journeys to destined host countries, including the most up-to-date situations and regulations on the way. This creates an opportunity for studying large-scale international migrations using the social media and data analytic methods, which have been

successfully applied in modelling various human behavior such as traveling tracking [2] and Crowd dynamics [3]. However, the migration pattern predication relies heavily on accurate identification of messages from refugees. This paper therefore focuses on identifying social media messages from refugees, which was confirmed as a challenging problem due to refugees reluctant to post on open social media platforms in the fear of getting caught by their government authority or human smugglers [4]. To prevent refugees' identities from being accidentally revealed by this study, their social media messages had been anonymized by replacing their usernames with unique index numbers.

Recent studies have analyzed the trends in mobility and migration flows using geolocated twitter data. For example, Zagheni [5] investigated the use of Twitter data to analyze the international and internal migration patterns. However, such analysis is entirely dependent on the availability of geolocated data, which is often unavailable in refugees' messages due to limitations on their computing devices or other reasons.

Gillespie [4] explored the refugee media journeys through smartphones and social media networks. It obtained information on trusted sources and groups on Facebook. While some of the Facebook groups studied contained news about refugees in Syria, other groups concerned general migration issues. These groups only revealed partial information about refugees' experiences. The study also found that most of the trusted sources communicates in Arabic.

Ali [6] discussed the significance of big data for various applications and development purposes. It provided a brief background of relevant techniques to understand the applications in humanitarian development. It proposed to use predictive analytics to avoid or mitigate humanitarian emergencies before they happened.

Brouckman and Wang [7] investigated the use of supervised machine learning along with Natural Language processing methods to classify downloaded tweets (with keyword 'refugee' in four languages) using the Twitter API. They used unigram features to model the tweets and then trained five classifiers namely, - Support Vector Machines, Logistic Regression, Random Forest, Naive Bayes, and Ensemble to predict the sentiments towards refugees as either positive, negative, or neutral.





