

Contents lists available at ScienceDirect

Data in Brief





Data Article

Medical dataset classification for Kurdish short text over social media



Ari M. Saeed^{a,*}, Shnya R. Hussein^a, Chro M. Ali^a, Tarik A. Rashid^b

ARTICLE INFO

Article history: Received 17 January 2022 Revised 14 March 2022 Accepted 16 March 2022 Available online 23 March 2022

Dataset link: Medical Sentiment Analysis Dataset for Kurdish Short Text over Social Media (Original data)

Keywords:
Machine learning
Medical text classification
Kurdish short text
Text pre-processing

ABSTRACT

The Facebook application is used as a resource for collecting the comments of this dataset, The dataset consists of 6756 comments to create a Medical Kurdish Dataset (MKD). The samples are comments of users, which are gathered from different posts of pages (Medical, News, Economy, Education, and Sport). Six steps as a preprocessing technique are performed on the raw dataset to clean and remove noise in the comments by replacing characters. The comments (short text) are labeled for positive class (medical comment) and negative class (non-medical comment) as text classification. The percentage ratio of the negative class is 55% while the positive class is 45%.

© 2022 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY license (http://creativecommons.org/licenses/by/4.0/)

^a Computer Science Department, University of Halabja, KRG, Halabja, Kurdistan, Iraq

^b Computer Science and Engineering Department, University of Kurdistan-Hawlêr, KRG, Erbil, Kurdistan, Iraq

^{*} Corresponding author.

E-mail address: ari.said@uoh.edu.iq (A.M. Saeed).

Specifications Table

Subject	Applied Machine Learning		
Specific subject area	Medical dataset classification for Kurdish short text over social media		
Type of data	Text		
	Figure		
	Table		
How the data were acquired	Facepager application is used for collecting the comments after configuring.		
Data format	Raw		
Description of data collection	Each post is separated accurately to describe the type of class (medical or		
	non-medical), then the link of the post is copied and pasted in the Facepager		
	application for gathering the specified comments.		
Data source location	Kurdish post link in Facebook Application		
Data accessibility +Repository name: Mendeley Data			
	Data identification number: 10.17632/f2yfz4r9fr.1		
	Direct link to the dataset: https://dx.doi.org/10.17632/f2yfz4r9fr.1		

Value of the Data

- This is an effort of collecting a dataset in the field of medical text classification for the Kurdish language. Moreover. It can be beneficial for supporting and modeling patient health systems, health policies, and regulations.
- The data is preprocessed and ready for implementation by those researchers and scholars who conduct research work on the Arabic Alphabet, such as Persian, Arabic, and Urdu.
- The dataset can be used with several preprocessing steps such as stemming and lemmatization.

1. Data Description

1.1. Data collection

In this era, the health of people is a serious subject that researchers work on it closely [5,6]. For this purpose, it is important to read humans' views over social media. In this work, the Facebook application is used as a social media for creating a proper MKD. Nevertheless, to say that for predicting the right sight of humans by using machines, a good resource (dataset) is necessary. As it is clear, there are so many channels, websites, and live posts that can be used for this purpose. The database in this work id consisted of 6756 samples, which are divided into two different classes (medical and non-medical). The samples were collected from various pages

Table 1Alphabet similarities among (Kurdish, Arabic, Persian, Urdu) Languages.

NO.	Kurdish alphabetic	Arabic Language	Persian Language	Urdu Language
1	ئ	1	i or i	i or i
2	ب	ب	ب	ب
3	پ	ث	Ų	پ
4	ث	ے	ت	ت
5	ح	ح	ث	ت
6	ত্	ζ	č	ث
7	ζ	خ	હ	ح
8	έ	ے		€
9	٤	ذ	ح خ	
10	ر)	٥	ح خ د
11	٦	ز	ذ	2
12	ز	<u>u</u>	ر	ż
13	ڗٛ	ů	ز	?
14	<u>u</u> u	ص	ۯ)
15	ů	<u>ض</u>	من	ر ژ
16		ط	ت ش	ز
17	ع غ	ظ	ص	ڗٛ
18	ف	۶	ت ض	u u
19	ڤ	خ خ	ط	ش
20	ق	ع غ ن	ظ	ص
21	ک	ق	۶	ے ض
22	گ	<u>ا</u>	ė.	ط
23	J	J	ع خ ف	ظ
24	J	م	ق	
25	م	ن	ک	ع غ ن
26	ن	٥	گ	ف
27	8	و	J	ق
28	و	ي	م	ن ک
29	ی	٠	ن	_ گ
30	١	•	و	J
31			ه	
32	و		ی	م ن
33	ۆ		3	
34				و ^
35	وو ى			~ ى
36	ێ			
טכ	G			۷

Table 2Number and percentage of collected dataset.

Class	Field	No. of Samples	Percentage
Medical	Medical	3076	45%
Not	News	890	55%
Medical	Economy	720	
	Education	1140	
	Sport	930	

and different areas as shown in Table 2. The number of medical comments (positive class) is 3076 while the non-medical comments (negative class) are 3680.

1.2. Methodology

On social media, the data can be viewed in various types, such as image, video, text. In this work, the data set is collected from the text. Facebook application is used for collecting the comments of users. Some different tools and techniques can be utilized for collecting the comments,

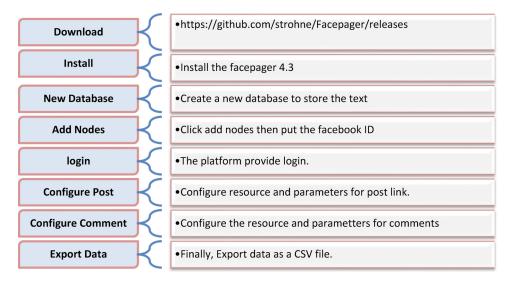


Fig. 1. Steps to dataset collection.

the Facepager tool is one of them that has been used for this reason [7]. The following steps should be followed for obtaining the data as shown below in Fig. 1.

As shown in Fig. 1, the first step is downloading the Facepager software for collecting the comments. The second step is locating and installing the files. The third step is to open the software and create a new database for saving the text file in (.db) format. The fourth step is adding nodes and putting the Facebook ID of the specified link after converting it over the internet. The fifth step is to log into Facebook via the Facepager tool. The sixth important step is configuring resources as (/<page-id>/posts) and parameters filed as (message) and specifying a start date and end date to fetch posts between those specific dates as shown in Fig. 2.

The seventh step is configuring a tool for fetching comments by clicking on a specific post and configuring resources as (/<post-id>/comments) and parameters filed as (message) as shown in Fig. 3.

The last and final step is exporting the comments as a CSV file as shown in Fig. 4.

1.3. Data set preprocessing

Preprocessing is one of the most important challenges for decreasing the noise on social media. Due to Kurdish users on the Facebook application using different Unicode to share their opinion and views. This causes a big issue for recognizing text and makes different characters shape. Using different scripts also increases the number of features (word) [1,4,8,9]. Accordingly, python language is used to create a new tool for implementing the below steps on the text as shown in Table 3:

- Removing noise (URL, User mentions, and Hashtag) on social media users will provide extra
 information for their relatives and friends by using URL, mentions (@user name), and hashtags (#special topic) that information are helpful for users but it is noise for the machine. It
 has to be removed.
- 2. Replacing elongated characters: users on social media sometimes use elongated words purposely to emphases about special things, such as (قوييويويون) (chiye), which means (Whaaaaaaat), which should be replaced with a base word (قوية) (chiye), which means (what).

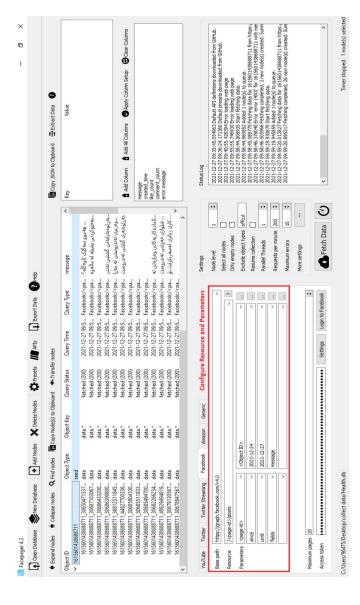


Fig. 2. Configure facepager posts.

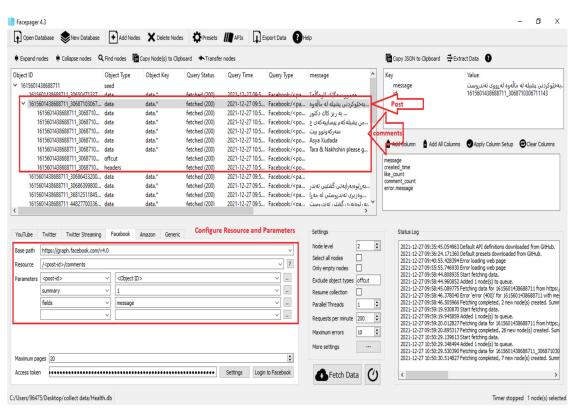


Fig. 3. Configure facepager comments.

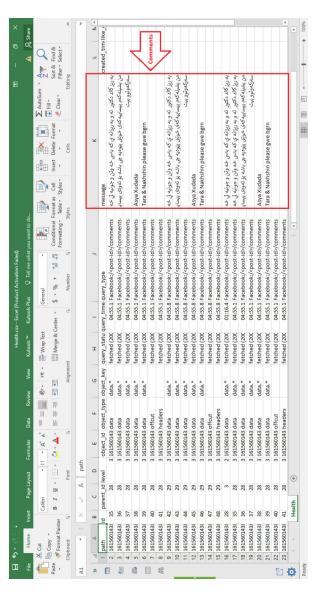


Fig. 4. Configure facepager comments.

Table 3 Preprocessing steps.

NO	Preprocessing steps	Natural comments with the Arabic alphabet	Natural comments with Latin alphabet	Natural comments in English	Preprocessing comments with the Arabic alphabet	Preprocessing comments with Latin alphabet	Preprocessing comments in English
1.	Removing noise (URL, User mentions and Hashtag):	داوز اوران پروتشد وروپسپ پرنهگدنپز و ورهگراستشراین واج ودلاهومگرفخن #Medical_Knowledge ولشروزپ ویتسرناز#	dktorej senariya ziyad zengene, psporî neşitergerî û nexoşiyekanî çaw. #Medical_Knowledge #zanistî_pzîskî	Doctor Sanaria Zyad Zangan. Expert in eye surgery and disease. Medical_Knowledge #نصريز پ	دایز ایرازاس اوروشکند عروبسرپ اینایگنزاوز و عربگرامششاین واج عناکهمشرفرخین	dktorej senariya ziyad zengenej psporî neşitergerî û nexoşiyekanî çaw	Doctor Sanaria Zyad Zangan expert in eye surgery and disease
2.	Replacing elongated characters	دې د مادې لکشروخ روتلکد والس توارګېب پر تس اسېداح تل پېځاېبلېب وودرده اتسي پښېلوجان پشراب کې پواج پېېېپېپېلګټ بېېپېپېېلکت بېښېپېواح دېپېې راګوه بېښېب د کېپېې	slaw dktur xuşkîkm heye le ĥadîsa serî bîkrawe aysta herdû bîlbîleyi cawî be başî naculînî tkayîyîyîyîyîyîyîye cwabibbibbibbib hukar cîyîyîyîye	Hi doctor, my sister was hurt in an accident now her cornea cannot move perfectly pleaaaaaaaaaase answeeeeer and whaaaaaat is the cause	هیده مکمیکشروخ رونتکد والس وارکمی، عررص اسیداح مل عیمای،بلی، وودرده انتسی، عنعلدوجان عشاب هب عواج دیج راکموه باوج دیاکت	slaw dktur xuşkîkm heye le ĥadîsa serî bîkrawe aysta herdû bîlbîleyi cawî be başî naculînî tkaye cwab hukar ciye	Hi doctor, my sister has hurt in an accident now her cornea cannot move perfectly please answer and what is the cause
3.	Incorrect spelling and grammars	و مین راج دُنْ مَچَرُدُهُ الْكُالْسُرام نوترزىھ بىساب مىارىباك تىخىبىد رىشنارگ شاكىب	maşalllla herçend car ew kabraye basî benzîn bikat grantir debêt	Allah has willed it, any time that man talks about gasoline, it will be more expensive	وهی راج دن هرده طل ااش ام ن عز ن هب ی س اب هی ارب اک ن کنب د رستن ارگ ساکب	maşaalle herçend car ew kabraye basî benzîn bikat grantir debêt	Allah has willed it, how any time that man talks about gasoline, it will be more expensive
4.	Removing punctuation	هتالُوو مهئ تتىبىغت ىھىقىىب !!!!مفىرىش ىماش	beqseyi tobêt em wullate şamî şerîfe!!!!	According to your speech, this country is peaceful!!!!	دتـــالْـــوو مەئ تــــئىبـــۇتـــ ىمســـق.دب دفـــــىردش ىماش	beqseyi tobêt em wullate şamî şerîfe	According to your speech, this country is peaceful
5.	Removing numbers:	یپوم طالس ۳۰ نم درونتگند دنستید ر ده مهنرید ر دد مهنریدر دو	dkturh min 30 salh mwî rdînm dh r dînm hh r dîth wh	Doctor, it is about 30 years, I have been pulling out beard hair, yet it grows back	دد م <i>ڼۍدر ټو</i> م طاس نم هروتکـد دو متــید ر ده مڼۍد ر	dkturh min salh mwî rdînm dh r dînm hh r dîth wh	Doctor, it is about 30 years, I have been pulling out beard hair, yet it grows back
6.	Replacing characters	عروز یمنځیساس نم والس عینتهبیات، کویده مواج راهیب یزرنوفل شیروځند وتوکتریبخدروس عینترفزیب عرضراچ ایای عیامزیور نای ځونه اینځمرفزیب	slaw min hsasiyetî z3rî çawm heye betaybetî lewerizî behar surdebîtewew dexurît aya çareserî bineretî heye yan rînmayî bîzehmet?	Hi, I have an eye rash, especially in the Spring season and it became red and itchy, is there any essential treatment or any advice.	کروز کیده کس اس نم والس کینسبی است وی ده مواج راهیب کرز و و ا کیب دروخه و موستی پیدروس کینسر وزیب کارسراج ای ای کیام کرکز رای دی ده	slaw min hsasiyetî zorî çawm heye betaybetî lewerizî behar surdebîtewew dexurît aya çareserî bineretî heye yan rînmayî bîzehmet	Hi, I have an eye rash, especially in the Spring season and it became red and itchy, is there any essential treatment or any advice.

- 3. Incorrect spelling and grammar: sometimes it is easy for users to correct the misspelling and grammar but machines cannot understand and it is challenging. These three words (الله الماله) (masha allh), which means (Allah has willed it) used as a misspelling instead the correct word (مالك القرام), which means (Allah has willed it).

- 5. Removing numbers: numbers increase the number of features in text datasets on social media and they are not helpful for the machine to understand. However. Kurdish users use different types of numbers, such as (English, Arabic, and Kurdish) numbers as shown.
- 6. Replacing characters: due to the Kurdish language using the same script of Arabic language for some characters and some users on social media use Arabic Keyboard for writing. This has become an issue for matching and selecting features. However, the issue has been solved by replacing the character as shown below:
 - a- 'ي' with 'ي'.
 - b- 'ك' with 'ك'.
 - c- 'ĕ' with 'ĕ'.
 - d- 'ö' with 'o'.
 - e- If the word ends with 'ن' replace with 'ن'.
 - f- If the word ends with ' \checkmark ' (\u 647\u 200C), then replace it with ' \checkmark ' (\u00b100647) as shown in the same shape of characters but different Unicode.

1.4. Dataset labeling

After collecting the dataset, another important step is labeling the samples. For this purpose, three annotators read the samples accurately and manually labeled the unlabeled samples for two classes (medical and non-medical). This process needs a huge effort and consumes time. For labeling each sample, the annotator annotates the sample based on some special words in the medical domain and the meaning of each sentence as shown in Table 4:

Table 4Labeling of comments.

NO.	Samples (Comments) in the Arabic alphabet	Samples (Comments) in Latin alphabet	Samples (Comments) in English	Classes
1	دی یہواق کالخ م واج وم دد نہ والس ت مح مزیب مشراب یج	slaw min de mu caw m xalî qaweyi ye cî başe bîze Äme t	Hello, my face has a brown spot, please what is good for me to do	medical
2	 آ یناک هب مل مت یهرامژ کـتولـونپ تیووم ده هزاو دئ تنب موهر مس زب اوین یـدووس چی.دهو زانام تنب 	polêk jmareyi te lh bih kanî ٦٠ bo sh rewe bê e wanh he muwî bê manan wehîç sûdî niye!	A class that has several students, more than 60, that is no sense and does't have any benefit	Not medical
3	ىكس دي دوام م ك ملانم روتكد والىس ى ريش دويب ى چنيل تتيج دئ تاوخ دئ وتتق	slaw dktur minalh kh m mawh yh skî ih cît lîncî biyuh şîr yi qtu ih xwat	Hi doctor, my baby has diarrhea and viscidity and eats condensed milk	medical
4	ى كىل كى كىلگىرەد ئىخىمب يىرقرۇر وەئ ھومئىاخەئ،ياد لىق،ئ مال،ىب ھومئىاك ھئ	ew rojeyi bext dergayt lê ekatewe belam eql dayexatewe	The day that opens the luck for you, yet, it closes mind	Not medical

Ethics Statement

All omments in the dataset belong to users in the Facebook application and it is scrapped. The data has been distributed over Facebook and thus, it has been collected and labeled. Moreover, we confirm that all the data is insensitive and anonymized data.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships which have, or could be perceived to have, influenced the work reported in this article.

Data Availability

Medical Sentiment Analysis Dataset for Kurdish Short Text over Social Media (Original data) (Mendeley Data).

CRediT Author Statement

Ari M. Saeed: Supervision, Data curation, Conceptualization, Methodology, Visualization, Project administration, Funding acquisition, Writing – original draft, Writing – review & editing; **Shnya R. Hussein:** Software, Formal analysis, Investigation, Resources; **Chro M. Ali:** Software, Formal analysis, Investigation, Resources; **Tarik A. Rashid:** Methodology, Supervision, Validation, Writing – review & editing.

Acknowledgments

The authors would like to thank the University of Halabja and the University of Kurdistan Hewler for providing all the facilities needed for conducting this research work.

References

- [1] A. Saeed, T. Rashid, A. Mustafa, R. Agha, A. Shamsaldin, N. Al-Salihi, An evaluation of Reber stemmer with longest match stemmer technique in Kurdish Sorani text classification, Iran J. Comput. Sci. 1 (2018) 99–107, doi:10.1007/s42044-018-0007-4.
- [2] T.A. Rashid, A.M. Mustafa, A. Saeed, A robust categorization system for Kurdish Sorani text documents, Inf. Technol. J. 16 (2016) 27–34, doi:10.3923/itj.2017.27.34.
- [3] T. Rashid, A. Mustafa, A. Saeed, Automatic Kurdish text classification using KDC 4007 dataset, in: advances in internetworking, data & web technologies. EIDWT 2017, Lecture Notes on Data Engineering and Communications Technologies, Springer, Cham, Wuhan, 2017.
- [4] A. Saeed, T. Rashid, A. Mustafa, P. Fattah, B. Ismael, Improving Kurdish web mining through tree data structure and Porter's Stemmer algorithms, UKH J. Sci. Eng. 2 (2018) 48–54, doi:10.25079/ukhjse.v2n1y2018.pp48-54.
- [5] R. Meena, V. hulasi Bai, Study on machine learning based social media and sentiment analysis for medical data applications, in: Proceedings of the Third International Conference on ISMAC (IoT in Social, Mobile, Analytics and Cloud) (ISMAC), IEEE, 2019.
- [6] G. Saranya, G. Geetha, C. K, M. Meenakshi K, S. Karpagaselvi, Sentiment analysis of healthcare tweets using SVM classifier, in: Proceedings of the International Conference on Power, Energy, Control and Transmission Systems (ICPECTS), IEEE, 2020.
- [7] W. He, S. Zha, L. Li, Social media competitive analysis and text mining: a case study in the pizza industry, Int. J. Inf. Manag. 33 (2013) 464–472, doi:10.1016/j.ijinfomgt.2013.01.001.
- [8] R. Ahmed, T. Rashid, P. Fatah, A. Alsadoon, S. Mirjalili, An extensive dataset of handwritten central Kurdish isolated characters, Data Brief 39 (2021) 107479, doi:10.1016/j.dib.2021.107479.
- [9] U. Naseem, I. Razzak, P. Eklund, A survey of pre-processing techniques to improve short-text quality: a case study on hate speech detection on twitter, Multimed. Tools Appl. 80 (2020) 35239–35266, doi:10.1007/s11042-020-10082-6.