

Building a corpus on Eating Disorders from TikTok: challenges and opportunities

Melissa Donati, Ludovica Polidori, Paola Vernillo and Gloria Gagliardi

Alma Mater Studiorum - University of Bologna, Italy

Abstract

We present two synchronic corpora of Eating Disorders (ED) related discourse on Social Media. PAC (i.e., ProAna/Anorexia Corpus) and RAC (i.e., Recovery from Ana/Anorexia Corpus) resources focus on the contents posted on TikTok, respectively, by communities promoting anorectic behavior and users sharing experiences concerning the process of recovery from their ED. We report on the corpus statistics and creation process, focusing specifically on the methodological issues raised by this novel Social Media platform.

Keywords

Eating Disorders, Corpus Linguistics, TikTok

1. Introduction

It was only 20 years ago that one of the darkest sides of Eating Disorders (ED) was revealed through the proliferation of websites, blogs, and social networks, in which a growing number of adolescents and young adults started sharing information about their eating experiences with like-minded users. Among these pro-ED communities, researchers and clinicians showed particular concern for pro-Ana (i.e., “pro-anorexia”) groups, i.e., web-based communities of anorexic (or aspiring anorexic) individuals engaged in the promotion of their Eating Disorder [1]. Interestingly, one of the most horrific and dangerous aspects of pro-Ana groups is that Anorexia Nervosa (AN) is not presented as a psychiatric disorder associated with pathological body image dissatisfaction [2], but more as a way of living with its own rules and rituals to be respected. While over the last years, much has been done to prevent the circulation of pro-ED content on social media (e.g., TikTok’s adoption of measures to obscure harmful contents: [3]), a new but specular phenomenon recently took the toll, that is, the spread of pro-recovery accounts of individuals who are in the process of healing from an ED and are willing to share their eating experience to help other online users [4]. From a linguistic perspective, research on ED has been very limited and became an object of study only in recent years [5, 6, 7, 8, 9, 10] as opposed to other psychopathologies, such as schizophrenia [11, 12], personality disorder [13], and depression [14, 15, 16, 17]. This already problematic picture has been

further compromised by the inhomogeneous representation of linguistic data in the literature, where the majority of studies have been dedicated to the linguistic profiling of ED-affected individuals in a Germanic language (English, German, Norwegian) [18]. This paper represents a small step towards the reversal of this tendency but a crucial part of two larger projects (Metaphan¹ and RaAM project 2022²) aiming at identifying, by the adoption of different NLP techniques and tools, potential lexical and semantic patterns in anorectic individuals. To this end, in the current research, we show the data collection process (i.e. oral and written productions) from ED communities on TikTok, currently representing the most widely used social media among young people and adolescents, namely the population groups at greater risk for EDs. In the following paragraphs, we give a brief overview of the literature on the topic (Section 2), then we describe the process of creating the corpus and discuss the methodological issues that were met (Section 3) and to conclude we provide few insights for future works (Section 4).

2. Related Works

In recent years, we have witnessed exponential growth in the use of Social Media (SM), especially by adolescents and young people. The community-building nature and the interactive dynamics of these platforms, as well as the less direct way of communicating, encourage users to openly discuss a wide variety of topics [19]. In turn, this makes available huge amount of data that can be used for different purposes (e.g. extract actionable patterns, form conclusions about users, conduct research, etc.). For this reason, Social Media Mining (SMM), i.e., the process of extracting big data from SM, now constitutes

¹<https://site.unibo.it/metaphan/en>

²<https://site.unibo.it/metaphan/en/connected-research-activities>

CLiC-it 2023: 9th Italian Conference on Computational Linguistics, Nov 30 – Dec 02, 2023, Venice, Italy

✉ melissa.donati@studio.unibo.it (M. Donati);

ludovica.polidori@studio.unibo.it (L. Polidori);

paola.vernillo@unibo.it (P. Vernillo); gloria.gagliardi@unibo.it

(G. Gagliardi)



© 2023 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

CEUR Workshop Proceedings (CEUR-WS.org)

a well-established methodology to collect large samples of data in different research areas [20]. This approach has proved particularly fruitful for collecting data on EDs as people suffering from these disorders seem to overcome the self-protective nature of their ED to engage in ED-related discourse with online users sharing similar experiences [21]. Indeed, in the last decade, many studies have used different SM platforms as a source of data to analyze EDs [22, 23, 24, 25, 21, 26, 27, 28, 29]. However, the state-of-the-art on ED-discourse on SM currently presents two main limitations: i) the majority of the analysis was carried out on small datasets built ad-hoc for the purpose of the work (with the only exception of [30]), and ii) they mostly focused on the English language. As a matter of fact, in the Italian framework there have been very little research on the representation of EDs on SM, and that little was mostly focused on Anorexia-Nervosa and did not target EDs in general [31, 32, 33].

3. Corpus Creation: Methodological Issues

Against this background and intending to fill this gap, we created a collection of English and Italian ED-related data that could be used for different types of research (from purely linguistic and content analyses that could help pinpointing the features and characteristics of ED-related discourse, to various computational techniques that could be used to implement systems of automatic detection of ED-related content on SM). We selected TikTok as a source of data as it currently represents the most widely used SM, especially among young people and adolescents, namely the at-risk population for EDs [34].

To achieve this goal, we first needed to define the nature and characteristics of the corpus itself. As far as the linguistic features are concerned, our corpus is specialized (i.e., is focused on the topic of EDs discourse on TikTok), synchronic (i.e., refers to a specific point in time that is the moment the data were downloaded), and targets both written and spoken language (TikTok videos contain spoken and/or written text). We did not set *a priori* a target dimension to be reached, because this feature is totally dependent upon the possibility of extracting the data automatically (Section 3.1). Conversely, following the common practice in the domain of SMM, we assumed that ‘there is no data like more data’ and intended to download as many videos as possible. To maximize the corpus representativity, we tried to balance the sample with respect to the types of videos being collected but we could not do so concerning the users’ gender, because for both corpora the vast majority of profiles were of female individuals (see Section 3.3 for more details). The target population consisted of those profiles that identify themselves in one of the two following categories: i)

supporters of anorectic behaviors (for the English PAC corpus); ii) witnesses and motivators for the recovery process (for the Italian RAC corpus). Such profiles were identified based on the linguistic and non-linguistic (i.e., emojis) information present in their profile bio. The selection criteria will be presented in Section 3.1, prior to the description of the data collection process and the discussion of the related issues that were encountered.

Before getting further into the methodology, it is necessary to make an ethical consideration concerning the collection of data from SM. Broadly speaking, SM posts that are publicly accessible are treated as belonging to the public domain, therefore, according to common practice, consent from the creators is not deemed necessary to download such data. This is strengthened by the fact that, upon registration, TikTok asks its users to consent to a set of terms of service that make the data available for access to third parties [35]. In addition, when creating and managing their accounts and contents users can decide to make them publicly accessible or private (i.e. only viewable by accepted followers); at any time, they can also restrict access to some of their contents through privacy settings and choose whether to make them downloadable. For the above reasons, given that for the purpose of this work only public and downloadable data was analysed, we did not seek users’ consent to collect the posts. In compliance with similar SM analysis [26], no reference to any identifying information, such as usernames, will be made.

3.1. Data Collection

As explained above, the selection criteria adopted to identify the target profiles was based on the information present in the profiles’ bio. However, to track the target profiles, we needed to start from a list of ED-related hashtags that could lead us to such profiles via a keyword-based search. The hashtags that were used herein were generated both by brainstorming and by exploring the platform for a couple of weeks, noting down the most popular trends and the most widely used hashtags (see Table 1 for an overview). Following this hashtag-driven search, we noticed that there was very little -if any- pro-Ana content produced in Italian, that is why for this type of ED-related content we decided to collect a small sample of English data. On the other hand, we found quite some profiles representing the ED-recovery community.

Among these profiles, we selected those having at least 10k followers (some of them exceed 2M followers) and at least 10 ED-related posts, so that we could maximize the chance of gathering interesting and relevant linguistic information. We then used the ED-related hashtags to conduct a within-profile research to select only the ED-related videos in each profile in order to extract them.

At this point, the next step consisted of extracting the

Table 1

List of pro-Ana and pro-Recovery hashtags that were used to search for TikTok profile that share ED-related content.

Pro-Ana hashtags	Pro-Recovery hashtags
#weightloss (w3ightl0ss)	#dcarecovery (dcar3covery)
#unhealthyweightloss (+ lexical variations)	#dca ⁴ #dcaitalia #fiocchettolilla ⁵
#kpop ³	#dcfighting

identified ED-related videos from the selected profiles. For the sake of time and efficiency, we wanted to download the data automatically. However, differently from other popular SM, TikTok has not yet released any official API that can be used by researchers and developers to automate the process of accessing and extracting the data. In addition, even if unofficial APIs exist, they get outdated almost immediately after their release because TikTok is constantly updating the anti-bot system preventing automatic access from the same IP. To get around this, we looked for a reliable and cost-effective proxy provider for TikTok scraping, but we could not find any viable solution.

Therefore we decided to proceed with the manual downloading of the data. The main drawback of this way of proceeding is that due to time and resource constraints we could not collect a very large number of videos (see Table 2). On the bright side, however, the manual downloading allowed us to i) enhance the content filtering process and ii) notice that TikTok videos have different formatting styles that might be worth distinguishing not only to ease the ensuing transcription process but also to conduct separate content analysis and compare the different results. Based on our observations about the different formatting styles, we grouped the TikTok videos into 4 subcorpora: 1) *Speech-only* videos: in which the user was talking in the absence of background music and/or written text; 2) *Playback*: in which the user lip-sync over a song or an extract from a movie or tv shows; 3) *Text-only*: in which there is neither background music nor the users themselves speaking, but only written text superposed on the video; and 4) *Mixed*: in which the above-mentioned features are present in various combinations.

³K-pop (for Korean-pop) is a popular genre of music originating from South Korea that has been hugely influential in the ‘diet scene’ because young people want to look like their favourite K-pop stars that are known for their extreme diets, indeed many young artists have left behind the K-pop world in order to focus on eating disorder treatment.

⁴Disturbo del Comportamento Alimentare (Eating Disorder).

⁵The Lilac Ribbon is the official international symbol against Eating Disorders.

3.2. Transcription

Organizing the videos into 4 categories was particularly useful for the transcription phase as it allowed to adopt different strategies and techniques based on the input characteristics. As for the downloading phase, although we intended to automatize the transcription process as much as possible, the high complexity of the data has, in some cases, made human intervention necessary.

For speech-only and playback videos automatic transcription was performed using the Google Web Speech API, which is easily accessible through the SpeechRecognition Library [36]. To assess the quality of the automatic transcription, a random sample of videos (n=10) for each category was extracted, transcribed manually and then compared with the machine-based transcription. For speech-only videos, a high agreement score was obtained between human and machine transcription (>90%) which confirmed the viability of the method adopted. Conversely, playback videos emerged as more problematic, thus manual correction was needed because both singing and the music accompaniment adversely impacted on intelligibility.

Automatic transcription was also attempted for text-only videos by means of Optical Character Recognition (OCR) using the Tesseract OCR engine [37], but we obtained poor results due to the high visual complexity of the input data, more specifically to the extreme variability of font type, size, and color, the lack of adequate contrast with the background, the non-hierarchical spatial organization of texts, and the presence of non-textual graphical elements (e.g., lexical variations of words, where letters are substituted by numbers or emojis to prevent the platform’s censorship and filtering system from blocking the content as potentially harmful, e.g., ‘starving’ written replacing star with the corresponding emojis, or ‘disorder’ written as ‘d1s0rder’). The same issue, boosted to the maximum, was observed with mixed videos, where speech, music, and written text were mingled. Therefore, for these two categories of videos, we could only perform the transcription manually.

We reported below, as an example of the type of ED-related content that was selected, the transcription of two videos, one for each of the two datasets.

[from RAC]

"questo video è davvero davvero difficile da registrare per me ma lo faccio perché voglio condire tutta la mia vita con voi e voglio aiutare delle persone che si trovano nella mia stessa situazione parlando del mio problema dovete sapere che io sono stata prima anoressica sono arrivata a pesare 36 kg e vi parlerò poi te la causa scatenante poi riscoperto il cibo ho iniziato ad abbuffarmi in una maniera assurda a sentirmi in colpa e quindi poi a vomitare questa

si chiama bulimia ovviamente alternavo momenti digiuno quindi magari non mangio proprio per giorni a momenti in cui il tuo corpo ha bisogno di cibo e quindi ti abbuffi e mangi qualsiasi cosa volevo solo dirvi che ieri è successa un'altra volta il fatto è che io me lo vedo subito in faccia cioè mi vedo 10 volte più grossa e mi sento davvero super gonfia che senti ma sono riuscita a non vomitare perché io sono più forte sono con tutte voi⁶

[from PAC]

*"i'm *** i'm a new member stats starter weight 140.1 ibs goal weight 100 ibs ultimate goal weight 90 ibs for now i binge eat when i'm bored so i gained a lot of weight in the past months i'm trying to limit myself on eating i am currently 4'10 and i'm overweight for my height age i listen to subliminal and trying to workout also i hate exercising but i realized it is healthy for me and my body 33"*

3.3. Corpus Statistics

In Table 2, we reported an overview of the statistics for the two corpora in terms of number of videos, number of words, and number of users from whose profiles the data were extracted.

The two corpora are registered in CLARIN⁷, but not publicly accessible for the moment.

Table 2
Statistics for the two corpora.

	PAC	RAC
n videos	250	1000
n words	13169	116261
n users	14 (all F)	27 (26 F, 1 M)

⁶[our translation] *"making this video is really really hard for me but I am doing it because I want to share everything about my life with you and I want to help those who are experiencing the same situation by talking about my problem you must know that I have suffered first from anorexia I ended up weighting 36 kg and I will tell you about the trigger then I rediscovered food and started insanely bingeing and feeling guilty and then as a consequence throwing up this is called bulimia obviously I alternated periods of fasting so peraphs I would not eat for days with periods in which my body needed food and I would eat anything and I just wanted to tell you that yesterday it happened again and the thing is that I see it immediately on my face that is I see myself 10 times bigger and I fell really extremely bloated that you know but I managed not to throw up because I am stronger I am with you all"*

⁷<http://hdl.handle.net/20.500.11752/OPEN-997>

4. Conclusion and Future Works

The aim of this work was twofold: on the one hand, we wanted to present two corpora on EDs, the English pro-Ana corpus (PAC) and the Italian pro-Recovery corpus (RAC), that were both built by extracting data from the popular SM TikTok; on the other, we wanted to discuss some methodological issues related to building a corpus using this platform as a source of data. More specifically, we pointed out that the absence of an official API does not allow the automatic extraction of the videos and requires manual work, which is highly time-consuming and does not allow to collect a very large sample of data. This, in turn, might impede the application of more complex computational analysis and limit the generalizability of the results. In addition, we raised the issue related to the transcription of the videos to text. In this case, implementing automatic approaches is not always feasible because of the extreme visual complexity and variability of TikTok videos.

Given the highly interactive nature of this SM and its unprecedented success, we believe that TikTok constitutes an extremely interesting source of linguistic and non-linguistic data that could be used to analyze other complex social and psychological phenomena and we hope that this work paves the way for further research in this direction.

CRedit authorship contribution statement

MD Conceptualization, Methodology, Software, Data Curation (i.e., download, automatic transcription, annotation), writing (§2,3,4)

LP Data Curation (i.e., manual transcription)

PV Conceptualization, Data Curation (i.e., download), Writing (§1)

GG Supervision, Funding acquisition.

Funding

This work was partially funded by the RaAM Association (project "How about metaphors for dinner? A digest of metaphorical conceptualizations in pro-Ana communities") and the University of Bologna (AlmaIdeas 2022 - "MetaphAN" project).

References

- [1] N. Boero, C. J. Pascoe, Pro-anorexia communities and online interaction: Bringing the pro-ana body online, *Body & Society* 18 (2012) 27–57.
- [2] J. B. Williams, M. First, Diagnostic and statistical manual of mental disorders, in: *Encyclopedia of social work*, 2013.
- [3] S. Marsh, Tiktok investigating videos promoting starvation and anorexia, *The Guardian* 7 (2020).
- [4] A. K. Greene, H. N. Norling, L. M. Brownstone, E. K. Maloul, C. Roe, S. Moody, Visions of recovery: a cross-diagnostic examination of eating disorder pro-recovery communities on tiktok, *Journal of Eating Disorders* 11 (2023) 109.
- [5] C. F. Bates, “i am a waste of breath, of space, of time” metaphors of self in a pro-anorexia group, *Qualitative Health Research* 25 (2015) 189–204.
- [6] O. Knapton, Pro-anorexia: Extensions of ingrained concepts, *Discourse & Society* 24 (2013) 461–477.
- [7] E. J. Lyons, M. R. Mehl, J. W. Penebaker, Pro-anorexics and recovering anorexics differ in their linguistic internet self-presentation, *Journal of psychosomatic research* 60 (2006) 253–256.
- [8] F. Skårderud, Eating one’s words, part ii: The embodied mind and reflective function in anorexia nervosa—theory, *European Eating Disorders Review: The Professional Journal of the Eating Disorders Association* 15 (2007) 243–252.
- [9] F. Skårderud, Eating one’s words, part i: ‘concretised metaphors’ and reflective function in anorexia nervosa—an interview study, *European Eating Disorders Review: The Professional Journal of the Eating Disorders Association* 15 (2007) 163–174.
- [10] M. Wolf, F. Theis, H. Kordy, Language use in eating disorder blogs: Psychological implications of social online activity, *Journal of Language and Social Psychology* 32 (2013) 212–226.
- [11] V. Bambini, G. Arcara, M. Bechi, M. Buonocore, R. Cavallaro, M. Bosia, The communicative impairment as a core feature of schizophrenia: Frequency of pragmatic deficit, cognitive substrates, and relation with quality of life, *Comprehensive psychiatry* 71 (2016) 106–120.
- [12] J. De Boer, M. Van Hoogdalem, R. Mandl, J. Brummelman, A. Voppel, M. Begemann, E. Van Dellen, F. Wijnen, I. Sommer, Language in schizophrenia: relation with diagnosis, symptomatology and white matter tracts, *npj Schizophrenia* 6 (2020) 10.
- [13] A. Arntz, L. D. Hawke, L. Bamelis, P. Spinhoven, M. L. Molendijk, Changes in natural language use as an indicator of psychotherapeutic change in personality disorders, *Behaviour research and therapy* 50 (2012) 191–202.
- [14] J. D. Bernard, J. L. Baddeley, B. F. Rodriguez, P. A. Burke, Depression, language, and affect: an examination of the influence of baseline depression and affect induction on language, *Journal of Language and Social Psychology* 35 (2016) 317–326.
- [15] T. Brockmeyer, J. Zimmermann, D. Kulesa, M. Hautzinger, H. Bents, H.-C. Friederich, W. Herzog, M. Backenstrass, Me, myself, and i: self-referent word use as an indicator of self-focused attention in relation to depression and anxiety, *Frontiers in psychology* 6 (2015) 1564.
- [16] N. Ramirez-Esparza, C. Chung, E. Kacewicz, J. Penebaker, The psychology of word use in depression forums in english and in spanish: Testing two text analytic approaches, in: *Proceedings of the international AAAI conference on web and social media*, volume 2, 2008, pp. 102–108.
- [17] J. Zimmermann, T. Brockmeyer, M. Hunn, H. Schauenburg, M. Wolf, First-person pronoun use in spoken language as a predictor of future depressive symptoms: Preliminary evidence from a clinical sample of depressed patients, *Clinical psychology & psychotherapy* 24 (2017) 384–391.
- [18] V. Cuteri, G. Minori, G. Gagliardi, F. Tamburini, E. Malaspina, P. Gualandi, F. Rossi, M. Moscano, V. Francia, A. Parmeggiani, Linguistic feature of anorexia nervosa: a prospective case–control pilot study, *Eating and Weight Disorders-Studies on Anorexia, Bulimia and Obesity* (2021) 1–9.
- [19] A. Lenhart, K. Purcell, A. Smith, K. Zickuhr, Social media & mobile internet use among teens and young adults. millennials., *Pew internet & American life project* (2010).
- [20] P. Gundecka, H. Liu, Mining social media: a brief introduction, *New directions in informatics, optimization, logistics, and production* (2012) 1–17.
- [21] T. E. Kenny, S. L. Boyle, S. P. Lewis, # recovery: Understanding recovery from the lens of recovery-focused blogs posted by individuals with lived experience, *International Journal of Eating Disorders* 53 (2020) 1234–1243.
- [22] M. Lukač, et al., Down to the bone: A corpus-based critical discourse analysis of pro-eating disorder blogs, *Jezikoslovlje* 12 (2011) 187–209.
- [23] L. Mullany, C. Smith, K. Harvey, S. Adolphs, ‘am i anorexic?’ weight, eating and discourses of the body in online adolescent health communication, *Communication & medicine* 12 (2016).
- [24] M. Moessner, J. Feldhege, M. Wolf, S. Bauer, Analyzing big data in social media: Text and network analyses of an eating disorder forum, *International Journal of Eating Disorders* 51 (2018) 656–667.
- [25] B. K. Bohrer, U. Foye, T. Jewell, Recovery as a process: Exploring definitions of recovery in the context of eating-disorder-related social media forums, *International Journal of Eating Disorders* 53

- (2020) 1219–1223.
- [26] S. S. Herrick, L. Hallward, L. R. Duncan, “this is just how i cope”: An inductive thematic analysis of eating disorder recovery content created and shared on tiktok using# edrecovery, *International journal of eating disorders* 54 (2021) 516–526.
 - [27] G. L. Jordan, M. D. Garcia, B. L. Diez, P. M. Sánchez, J. G. Del Barrio, R. Ayesa-Arriola, Facebook as a pro-ana and pro-mia resource, *European Psychiatry* 64 (2021) S703–S703.
 - [28] C. González-Nuevo, M. Cuesta, J. Muñiz, Concern about appearance on instagram and facebook: Measurement and links with eating disorders, *Cyberpsychology: Journal of Psychosocial Research on Cyberspace* 15 (2021).
 - [29] M. Minadeo, L. Pope, Weight-normative messaging predominates on tiktok—a qualitative content analysis, *Plos one* 17 (2022) e0267997.
 - [30] M. Donati, C. Strapparava, CorEDs: A corpus on eating disorders, in: *Proceedings of the RaPID Workshop - Resources and Processing of linguistic, para-linguistic and extra-linguistic Data from people with various forms of cognitive/psychiatric/developmental impairments - within the 13th Language Resources and Evaluation Conference, European Language Resources Association, Marseille, France, 2022*, pp. 80–85. URL: <https://aclanthology.org/2022.rapid-1.10>.
 - [31] V. Richichi, A. Chinello, F. Parma, L. E. Zappa, E. Mazzoni, F. Monti, Anoressia nervosa e internet. uno studio sui blog pro-ana in italia, *Psicologia clinica dello sviluppo* 22 (2018) 499–514.
 - [32] N. L. Bragazzi, G. Prasso, T. S. Re, R. Zerbetto, G. Del Puente, A reliability and content analysis of italian language anorexia nervosa-related websites, *Risk management and healthcare policy* (2019) 145–151.
 - [33] G. Gagliardi, “odio tutto ciò, voglio le ossa”: Una prima indagine sulle caratteristiche linguistiche delle pagine social pro-ana in lingua italiana, *Italiano LinguaDue* 13 (2021) 520–536.
 - [34] A. Sherman, Tiktok reveals detailed user numbers for the first time, Retrieved October 2 (2020) 2020.
 - [35] 2023. URL: <https://www.tiktok.com/legal/page/eea/privacy-policy/en>.
 - [36] Y. Zhang, M. Pezeshki, P. Brakel, S. Zhang, C. L. Y. Bengio, A. Courville, Towards end-to-end speech recognition with deep convolutional neural networks, *arXiv preprint arXiv:1701.02720* (2017).
 - [37] J. Ooms, tesseract: Open Source OCR Engine, 2023. <https://docs.ropensci.org/tesseract/> (website) <https://github.com/ropensci/tesseract> (devel).