

ソーシャルメディアテキストを用いた摂食障害の文化差比較

栗生 紗希帆¹ Daveon Kim² Hvuk-Yoon Kwon² 若宮 翔子¹ 荒牧 英治¹

¹ 奈良先端科学技術大学院大学 ² Seoul National University of Science and Technology

{kuriu.sakiho.ks0, wakamiya, aramaki}@is.naist.jp

{rlaekdus7668@ds., hyukyoan.kwon}@seoultech.ac.kr

概要

日本と韓国の肥満率は、先進国の中でワースト1位、2位である。過度なやせは摂食障害をもたらすことがある。摂食障害は、身体面、心理面、行動面への治療が必要であり、数ヶ月から数年の時間を要するといわれている。本研究では、日本と韓国を対象に、摂食障害とダイエットのソーシャルメディアテキストを収集し、分類した。そして、摂食障害やダイエットに関連する言語的特徴は、文化間（日本と韓国）でどのように異なるのかを調査した。結果としては、日本と韓国による言語的な文化差がみられ、2言語のテキストを活用した方が性能が僅かに向上した。

1 はじめに

近年、日本および韓国においては、アイドル文化が顕著な経済的影響を及ぼしている。2022 年における日本の経済効果は 1,650 億円とされ、2023 年には 2,000 億円近くに拡大することが予想されている [1]。さらに、2017 年から 2021 年にかけての世界的な「韓流」ブームを背景に、韓国の輸出が増加し、これに伴う韓国の経済効果は 37 兆ウォン（約 4 兆 380 億円）に達するとの報告がある [2]。しかしその一方で、外見に対する基準は極めて厳格であり、体型や体重を適切に管理できないことはしばしば努力不足として非難される。このような傾向は、日韓合同ガールズグループやアイドルオーディション番組の登場によって一層顕著になり、競争が激化している。結果として、極端な食事制限や過度な運動など、健康を害するリスクを伴うダイエットが行われる事象が増加している。さらに、近年ではインターネットの普及により、情報へのアクセスが一層容易になった結果、アイドルのダイエット方法が広範囲にわたり拡散しやすくなった。過酷なダイエットは、必要な食事量を大きく下回ることがあり、**摂食**

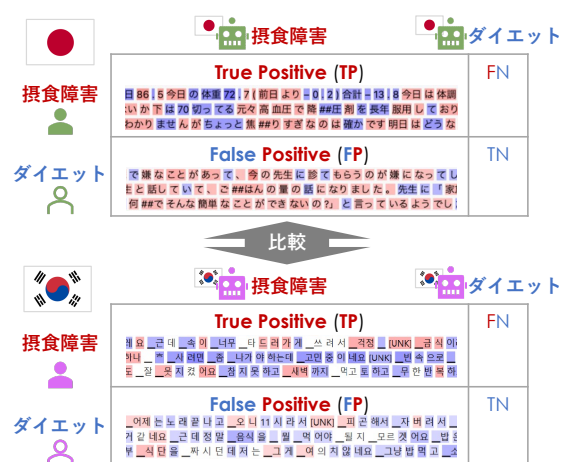


図1 本研究の概要。日本と韓国を対象に、摂食障害とダイエットに関連するソーシャルメディアテキストを収集、分類した結果から、文化間による言語的特徴の差異を議論する。

障害を引き起こすリスクが懸念される。

摂食障害とは、食行動の異常がみられ、精神的・身体的に影響が及ぶ病気である。さらに、摂食障害は**神経性やせ症**と**神経性過食症**に分けられる。神経性やせ症とは、食べることを極端に制限し、体重が減っても増加を恐れ、低体重を維持しようとする病気である。神経性過食症とは、食のコントロールができず、頻繁に過食し、その後後悔して嘔吐や下剤を使い、体重増加を防ごうとする病気である。日本国内の患者数は約 22 万人 [3] と推定されている。さらに、韓国では 2 万人を超えており、直近 5 年間の摂食障害患者数が 30～60 % 増加している [4]。摂食障害の治療には、軽度の場合には半年から 1 年程度で改善することもあるが、中等度から重度の場合は数年かかることがある。このように、摂食障害のリスクは高く、治療に時間がかかるため、その兆候を早期発見し、効果的な介入が極めて重要である。

本研究では、日本と韓国を対象に、摂食障害者によるテキストとダイエットに関連するテキストを収

表1 日本語と韓国語のデータセットの内訳。データ数は、日本語データセットではブログを、韓国語データセットではスレッドを示す。

| 言語 | データ種別 | データ収集期間 | データ投稿期間 | データ数 | ユーザ数 | データソース |
|-----|-------|--------------|----------------------|---------|--------|------------|
| 日本語 | 摂食障害 | 2023/6～8 | 2005/6/4～2023/8/1 | 81,340 | 380 | アメーバブログ |
| | ダイエット | 2024/8/3～8/5 | 2001/1/1～2024/8/5 | 407,668 | 826 | |
| 韓国語 | 摂食障害 | 2024/9～10 | 2009/6/13～2024/10/10 | 29,789 | 9,098 | Naver Cafe |
| | ダイエット | 2024/8～10 | 2004/2/6～2024/10/28 | 45,793 | 11,756 | |

集し、これらを用いて、摂食障害とダイエットに関連するテキストを分類する。さらに、単言語モデルと多言語モデルの両方を活用することで、それぞれのモデルの特性を生かし、日本と韓国の文化間における摂食障害やダイエットに関連する言語的特徴の違いを明らかにすることを目的とする。

2 関連研究

摂食障害の発症に関する要因は、社会文化的・家族的・心理的な側面に多岐にわたることが分かっている [5] が、研究としては、遺伝的要因や精神疾患、そしてボディイメージの歪みが中心課題となっている [6]。特に、日本と韓国における摂食障害の原因についての研究では、「ボディイメージ」に起因している可能性が示唆されている [7]。我々は、この自己のボディイメージに対する否定的な認識に注目している。

一方で、摂食障害の早期発見において、自然言語処理技術は注目を集めている [8]。ただし、英語データが研究の中心であり、異なる文化圏への対応が課題となっている。患者が残したテキストメッセージを用いた研究もあるが [9, 10]、これらも主に英語を対象とするにとどまっている。

本研究では、瘦身文化が根強い日本と韓国にて、ボディイメージ（ダイエット）における摂食障害のリスクにアプローチする。

3 データセット

本研究のデータセットは、日本と韓国の2カ国における、摂食障害とダイエットに関するタイトル、テキスト、投稿日時で構成されている。データセットの内訳を表1に示す。なお、収集したテキストへの前処理として、顔文字や全角スペースなどの特殊文字や、HTML タグなどを削除した。

日本語データセットでは、闘病記の「摂食障害」

表2 パラメータ設定

| パラメータ | 単言語モデル | | 多言語モデル |
|----------------|----------------------|----------------------|----------------------|
| | 日本語 | 韓国語 | |
| Max Length | 128 | 128 | 128 |
| Batch Size | 16 | 16 | 16 |
| Epochs | 100 | 20 | 100 |
| Learning Rate | 2.0×10^{-5} | 1.0×10^{-5} | 2.0×10^{-5} |
| Early Stopping | 2 patience | 3 patience | 3 patience |
| Optimizer | AdamW | | |
| Loss Function | Cross-entropy Loss | | |

「拒食症」¹⁾を摂食障害データとした。アメーバブログの公式ジャンル「ダイエット記録」²⁾をダイエットデータとして収集した。

韓国語データセットでは、Naver Cafe から摂食障害データ³⁾とダイエットデータ⁴⁾をそれぞれ収集した。Naver Cafe は、韓国最大の検索エンジンサービスを提供する「NAVER」が運営するコミュニティサイトであり、幅広い分野にわたる情報交換が可能である。あるユーザが、カフェ（コミュニティ）を開設して掲示板を作成し、カフェのユーザはその掲示板で自由にコミュニケーションを行う仕組みとなっている。

4 実験

4.1 設定

本研究では、摂食障害テキストとダイエットテキストを分類する。その後、分類モデルがどの特徴量に寄与しているかを明らかにするため、Integrated Gradients を用いて分析した。Integrated Gradients は、予測における重要な特徴量を特定する手法である。このとき、False Positive (FP) と True Positive (TP) のテキストに着目する。これは、言語モデルが「摂食

1) 闘病記「摂食障害」「拒食症」<https://www.tobyo.jp/>

2) アメーバブログ公式ジャンル「ダイエット記録」
<https://blogger.ameba.jp/genres/diet/blogs/ranking>

3) Naver Cafe <https://cafe.naver.com/jahayun>

4) Naver Cafe <https://cafe.naver.com/dietfood>

障害」に分類したテキストについて、どの特徴量が寄与しているかを明らかにし、日本語と韓国語それぞれの摂食障害に関する言語的特徴を議論するためである。データには、摂食障害またはダイエットのラベルが付与されている。モデル学習時のハイパーパラメータは、それぞれ表 2 のように設定した。分類モデルの評価指標には、Accuracy, F1-score, Recall, Precision を用いた。

4.2 単言語モデル

日本語および韓国語をそれぞれの言語に特化したモデルで訓練し、分類性能を検証した。

日本語分類モデルは、東北大学の事前学習済み BERT モデルである `cl-tohoku/bert-base-japanese-v3`⁵⁾ を用いた。日本語データセットからランダムに抽出した 137,500 件を用い、5 分割交差検証した。韓国語分類モデルは、韓国語の事前学習済み BERT モデルである `skt/kobert-base-v1`⁶⁾ を用いた。韓国語データセットからランダムに抽出した 50,045 件を用い、4 分割交差検証した（分割数が日本語に比べて少ないのはデータサイズを考慮してのことである）。

4.3 多言語モデル

日本語および韓国語を合わせて学習することで、性能が向上するかどうかを検証した。まず、韓国語データを用いてモデルを訓練した。その後、日本語データで再訓練した。さらに、比較実験として、日本語データ単体でもモデルを訓練した。この手法により、片方の言語がもう片方に貢献するかを検証できる。

分類モデルは、多言語事前学習済みモデルである `xlm-roberta-base`⁷⁾ を基に構築した。使用データは、日本語データセットからランダムに取り出した 58,000 件を用い、5 分割交差検証した。一方、韓国語データセットでは、日本語の訓練データ数と一致させるため、ランダムに 46,000 件を抽出して多言語モデルの訓練に用いた。

5 結果と考察

5.1 単言語モデルでの分類

分類モデルの性能を表 3 に示す。日本語モデル、韓国語モデルともに Accuracy, F1-score, Recall およ

表 3 単言語モデルの分類性能

| モデル | Accuracy | F1-score | Recall | Precision |
|-----|----------|----------|--------|-----------|
| 日本語 | 0.812 | 0.809 | 0.812 | 0.831 |
| 韓国語 | 0.845 | 0.846 | 0.845 | 0.845 |

表 4 多言語モデルの分類性能。「日本語」は日本語データのみで訓練したモデル、「韓→日」は韓国語データで訓練した後、日本語データで追加訓練したモデルを示す。

| モデル | Accuracy | F1-score | Recall | Precision |
|-----|----------|----------|--------|-----------|
| 日本語 | 0.990 | 0.990 | 0.990 | 0.990 |
| 韓→日 | 0.991 | 0.991 | 0.991 | 0.991 |

び Precision のスコアが 80% を超えており、高水準の性能を示している。韓国語モデルでは Precision と Recall がほぼ一致している一方で、日本語モデルでは Precision が優位となり、False Negative を多く検出している。

次に、Integrated Gradients を適用した FP および TP の 2 例を図 2 に示す。なお、韓国語のテキストの日本語訳については、図 3 に示す。ここで、日本語の分類に大きく寄与しているトークンは以下であった。FP に関連する特徴量には、「バイト」「ごはん」「診」「病気」「呼吸」「起き」「食」「悲」「抜け出し」「体」「夕飯」「買い物」「痛」などがある。TP に関連する特徴量には、「食」「呼吸」「吐き」「頃」「気温」「予想」「一番」「病気」「勉強」「参加」「面倒」「怖」「体重」「病院」「制限」などが挙げられる。FP においては、「バイト」「ごはん」「診」「買い物」などの語は、摂食障害を示唆するものでない用語を重視している。今後、データを十分に増やすことで、このような不適切な単語への重みづけが改善されと考えている。一方で、TP の事例では、「吐き」「体重」「制限」などのトークンを手掛かりにできている。韓国語の分類には、FP に関連する特徴量として、「ストレス」や否定形などの特徴がみられた。一方で、TP に関連する特徴量としては、「断食」「治療」など摂食障害の特徴的な用語が大きく寄与したことが読みとれた。

5.2 多言語モデルでの分類

分類モデルの性能を表 4 に示す。韓→日は、韓国語モデルを基に日本語データで再訓練したモデルであり、日本語のみの性能よりも僅かに高い。このことから、韓国語の知識を一部活用しながら日本語データに適応できる可能性がある。今後、他の文化、言語圏での実験を予定している。

5) <https://github.com/cl-tohoku/bert-japanese>
6) <https://github.com/SKTBrain/KoBERT>
7) <https://huggingface.co/FacebookAI/xlm-roberta-base>

JPJ012425 および花王株式会社共同研究費の支援を受けたものである。

参考文献

- [1] 株式会社矢野経済研究所. 「オタク」市場に関する調査を実施 (2023 年), 2023. https://www.yano.co.jp/press-release/show/press_id/3383.
- [2] 한국경제인협회 경제조사팀(韓國經濟人協會經濟調査チーム). 한류 확산의 경제적 효과(韓流の普及による經濟効果), 2023. https://www.fki.or.kr/main/news/statement_detail.do?bbs_id=00035093&category=ST.
- [3] 厚生労働省. 精神保健福祉資料：2021 年 ndb データ, 2021. <https://www.ncnp.go.jp/nimh/seisaku/data/ndb.html>.
- [4] 국민건강보험공단(國民健康保險公團). 국민건강보험공단_식이장애 관련 상병별 진료현황_20221231 (國民健康保險公團_摂食障害関連疾患別診療状況), 2024. <https://www.data.go.kr/data/15127691/fileData.do>.
- [5] Janet Polivy and C Peter Herman. Causes of eating disorders. **Annual review of psychology**, Vol. 53, No. 1, pp. 187–213, 2002.
- [6] Ruth H Striegel-Moore and Cynthia M Bulik. Risk factors for eating disorders. **American psychologist**, Vol. 62, No. 3, p. 181, 2007.
- [7] 野上真央, 吉村英一, 鈴木真理子, 田尻絵里, 中下千尋, 濱田有香, 塩瀬圭佑, 阿曾 (染矢) 菜美, 畑本陽一, 田中亮ほか. 若年やせ女性が形成される要因に関するスコーピングレビュー. 女性心身医学, Vol. 29, No. 2, pp. 206–219, 2024.
- [8] Ghofrane Merhbene, Alexandre Puttick, and Mascha Kurpicz-Briki. Investigating machine learning and natural language processing techniques applied for detecting eating disorders: a systematic literature review. **Frontiers in Psychiatry**, Vol. 15, p. 1319522, 2024.
- [9] Miguel Rujas, Beatriz Merino-Barbancho, Peña Arroyo, and Giuseppe Fico. Development of a natural language processing-based system for characterizing eating disorders. In **IberLEF@ SEPLN**, 2023.
- [10] Stella Maćkowska, Bartosz Koścień, Michał Wójcik, Katarzyna Rojewska, and Dominik Spinczyk. Using natural language processing for a computer-aided rapid assessment of the human condition in terms of anorexia nervosa. **Applied Sciences**, Vol. 14, No. 8, p. 3367, 2024.

A 付録

어제 는 노래 끝 나고 오 니 11 시 라 서 [UNK] 피곤 해서 자 버 려 서 쓰 지를 못 했 네요 [UNK] 어제 는 그렇게 폭 식 도 안 했 고 거의 노래 만 4 시간 한 거 같 네요 근 데 정말 음식 을 뭘 먹 어야 될 지 모르 겠 어요 밥 은 반 으로 줄 이 긴 헛 는데 나머 지 는 어떻게 해야 될 지 몰 라 서 다른 분 들은 전 부 식 단 을 짜 시 던 데 저 는 그 게 여 의 치 않 네요 그냥 밥 먹 고 소화 다 했 을 때 즈 에 운동 을 하는데 제 가 땀 이 많은 체 질 이라 40 분 정 도 하면 옷 을 갈 아 입 어야 될 정 도 [UNK] 근 데 노래 하는데 칼 로 리 소 모 가 많 을 까요 궁금 하 네요 분명 힘들 긴 한 데 얼마 나 소 모 가 될 지 아 시 는 분 계 신 가 요 [UNK] 그리고 마 칭 가 님 말씀 대로 저녁 6 시 전에 밥 반 공 기 라도 먹 어야 겠 네요 안 그 림 더 안 좋 을 것 같 아 요 [SEP] [CLS]

昨日は歌が終わって帰ったら11時になって 疲れて寝てしまったので書けませんでした [UNK] 昨日は爆食もせずほとんど歌だけで4時間過ごした気がします でも何を食べるべきか本当にわかりません 半分に減らしたけど 残りはどうすべきかがわからなくて 他の人たちは皆食事プランを作っているのに私はそれがうまくいきません ただご飯を食べて消化が全て終わった頃に運動をするんですが 私は汗をかきやすい体質なので 40分くらいすると服を着替える必要はないほど [UNK] でも歌っているときにカロリー消費が多いのでしょうか 気になりますね 確かに疲れますが どれくらい消費されるのかわっている人がいるのでしょうか [UNK] そしてマジンガさんの言う通り 夕方6時前にご飯半分でも食べなければなりませんね そうしないともっと良くないと思います [SEP] [CLS]

(a) 正解ラベル：ダイエット，予測ラベル：摂食障害

내 일 일 플 란 트 4 개 수 면 치료 받 으 려 가 네요 현 두 군 두 군 아침 10 시부터 가 님 폭 자 고 인 타 겠 어 [UNK] 현 눈 뜨 면 치 아 가 느껴 질 생 각 에 설 레 요 근 데 속 이 너무 타 드 러 가 게 쓰 러 서 걱정 [UNK] 금 식 이라 아침 못 먹 으 니 저녁 먹 어야 하는데 지금 딱 명 기 는 거 리 곤 커 스타 드 크림 빵 요 거 하나 현 사 러 면 좀 나가 야 하는데 고 민 중 이 네요 [UNK] 빈 속 으로 가 서 치료 끝 내 면 기 운 빠 저 서 백 프로 빈 혈 날 텐 데 [UNK] 전 연 누구 와의 약속 도 잘 못 지 켜 어요 참 지 못 하고 새벽 까지 먹고 토 하고 무 한 반 복 하고 약속 시간 몇 시간 전 인 데 도 폭 식 하고 토 하고 그러 다 보니 잠 못 자 고 피곤 해서 약속 파 토 내 고 지 각 하고 오늘 은 무 사 히 내 일 치료 만 생각 하며 바쁘 나 가 겠 습 니다 으 아 찻 [SEP] [CLS]

明日 インプラント4本 睡眠 治療で受けに行きます ドキドキしますが朝10時から行くので しっかり寝て起きなければいけません [UNK] 目が覚めたら歯が整っているのを感じると思うワクワクします でも胃がとてもむかむかして痛くて心配 [UNK] 断食なので朝食を食べられないので夕食を食べなければならないので すが今食べたいのはカスタードクリームパンだけ 買に行かなければならないので悩んでいます [UNK] 空腹時に行って治療が終わったら気力が落ちて間違いなく貧血になりそうです [UNK] 以前は誰かとの約束もよく守れませんでした 我慢できずに夜中まで食べて吐いてを無限に繰り返し 約束の時間 数時間前なのに過食して吐いて そうしていたら眠れずに 疲れて約束をすっぽかしたり遅刻したり 今日は無事に明日の治療だけを考えながら切り替えていきます よしっ [SEP] [CLS]

(b) 正解ラベル：摂食障害，予測ラベル：摂食障害

図3 韓国語モデルにおける Integrated Gradients の可視化結果の日本語訳。(a)は False Positive に、(b)は True Positive に分類されたテキストデータ。各行において、上段に韓国語、下段に日本語訳を示す。