

**2.【研究計画】** ※適宜概念図を用いるなどして、わかりやすく記入してください。なお、本項目は1頁に収めてください。様式の変更・追加は不可。

### (1) 研究の位置づけ

特別研究員として取り組む研究の位置づけについて、当該分野の状況や課題等の背景、並びに本研究計画の着想に至った経緯も含めて記入してください。

#### 当該分野の状況: フェイクニュースの自動検出

SNSの発展で情報を迅速かつ大量に取得・共有が容易になった一方、悪意により他人を騙すために作られたフェイクニュースも拡散されやすくなった。特に2020年からCOVID-19の影響による誤情報の拡散であるインフォデミックにより、メタノール飲用による死亡事故[1]といった事象が報告された。以上から騙された人々により社会的損害が起きるため、フェイクニュース拡散の早期抑制が必要である[2]。

フェイクニュース検出へ有識者が調査するファクトチェックがある。これは拡散ののち着手されるため、拡散抑制にはならない。そのため、自動でフェイクニュースを検出するべく深層学習によってファクトチェック結果をラベルとして、記事内容やユーザの反応から教師あり学習で自動検出する研究がある[3]。

#### 課題

フェイクニュース自動検出が抱える課題は以下の通りである。

##### 早期検出と正確性の両立

記事内容に加えてユーザの反応(RT, コメント等)を扱うと検出性能が改善した報告がある[4]一方、ユーザの反応を十分に得るには時間がかかるため、高い正確性と早期検出を両立できない。

##### 日本語データセット不足

深層学習による実現は、正解ラベルとして多量のファクトチェック結果を要する。このファクトチェックが活発な地域差の影響でデータセットが英語に集中[5]している。もし日本語を対象とした場合、ファクトチェック結果が不足しているためラベル付きデータセットによる教師あり学習ができない。

##### データセットのタイムリー性

ニュースという性質上、時間経過によって扱われる情報や出来事が徐々に古くなり、フェイクニュースの傾向変化への対応が徐々に難しくなる。実際に学習・検証で入力されたニュースが扱う出来事の違いによって検出性能の劣化が指摘される[3]。よってデータセットとモデルの提供のみならず、継続してデータセットを拡大する仕組みも必要である。

#### 本研究計画の着想に至った経緯

私は修士過程で英文フェイクニュース早期検出の研究を行った。記事に対するユーザのコメントが検出に有用とする先行研究をベースに、早期検出を想定して少ないコメントから更にコメント内容を自動生成して検出するモデルを実装した。実験にてコメントを生成した上で分類することでより多くのフェイクニュース検出を実現した(査読付き海外IEEE学会発表済[7])。

一方、国内研究会で発表したところ想定以上に日本語での実現に対する期待を受けた。日本は英語圏に比べファクトチェックされた記事が少なく、データセットを作ってモデルを実装するにはラベルが足りない。このラベル不足を補う方法として、少量のラベル付き記事と多量のラベルなし記事にユーザの初期反応から弱いアノテーションを付加して学習を行う弱教師あり学習を行う研究[8]に着目した。日本語で同じ構成のデータセットを作成し、分類を行うモデルを実装することで実現可能と考えた。

参考文献

- [1] H H-M, et al. *Critical Care* 24.1 2020: 1-3.  
 [2] S. Tasnim, et al. *JPMMPH* 53.3 2020: 171-174.  
 [3] Yaqing W, et al. *KDD'18*, pp. 849-857. 2018.  
 [4] Liang W & Huan L. *WSDM '18*, pp. 637-645, 2018.

- [5] Kai S, et al. *Big Data* 8.3 2020: 171-188.  
 [6] Yejin B, et al. *arXiv preprint arXiv:2101.03841* 2021.  
 [7] Yuta Y, et al. *INES*. 2020.  
 [8] Kai S, et al. *arXiv preprint arXiv:2004.01732* 2020.

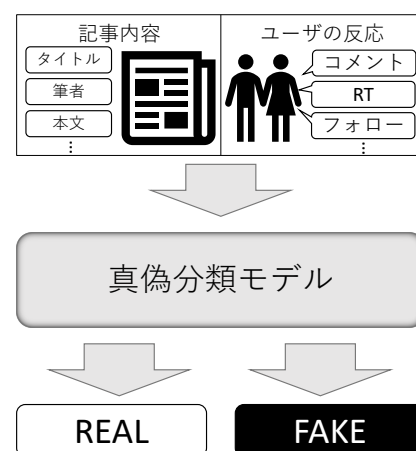


図1: フェイクニュース自動検出の基本的な流れ

【研究計画】（続き） ※適宜概念図を用いるなどして、わかりやすく記入してください。なお、各事項の字数制限はありませんが、全体で2頁に収めてください。様式の変更・追加は不可。

## (2) 研究目的・内容等

- ① 特別研究員として取り組む研究計画における研究目的、研究方法、研究内容について記入してください。
- ② どのような計画で、何を、どこまで明らかにしようとするのか、具体的に記入してください。
- ③ 研究の特色・独創的な点（先行研究等との比較、本研究の完成時に予想されるインパクト、将来の見通し等）にも触れて記入してください。
- ④ 研究計画が所属研究室としての研究活動の一部と位置づけられる場合は申請者が担当する部分を明らかにしてください。
- ⑤ 研究計画の期間中に受入研究機関と異なる研究機関（外国の研究機関等を含む。）において研究に従事することも計画している場合は、具体的に記入してください。

### ①研究計画における研究目的、研究方法、研究内容

#### 研究目的

本研究では、フェイクニュースの早期自動検出を日本語で実現・研究の促進を目指し、データセットの作成から検出モデルの実装を目的とする。さらにユーザに直接反応を伺うシステムを構築することで、早期検出の実現とデータセットの規模拡大と改善を続ける。

#### 研究方法・研究内容

以下の3目標を目指し研究する。

**目標Ⅰ** データセット作成に向けファクトチェック済記事とそうでない記事、同時に SNS 上で記事に寄せられたコメントやユーザ情報等を収集する。

**目標Ⅱ** 早期検出を想定した状況で高精度な真偽分類を行うモデルをユーザの初期反応から得られた弱いアノテーションと共に行う。

**目標Ⅲ** 新たな検出の枠組みとして利用者とのやり取りによってフェイクニュースの疑いがある記事の検出を行うモデルを開発する。

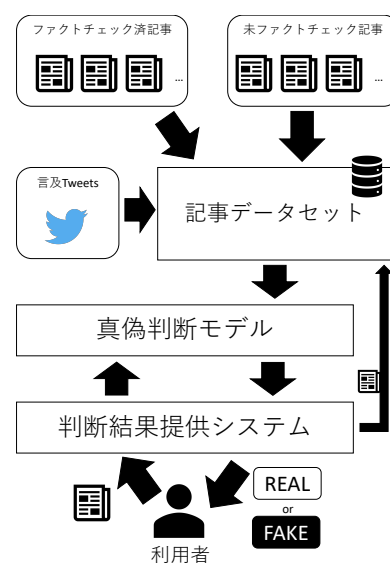


図2: 本研究計画の完成予想図。

### ②どのような計画で、何を、どこまで明らかにしようとするのか

**目標Ⅰ: 日本語の記事・真偽を含むデータセットを作成する**

(採用前 - 1年目)

日本語での検出を実現するためには、まずはデータセットを作成する必要があります。データセット作成の全体の流れは図2の通りである。日本語ファクトチェック結果の取得には、特定非営利活動法人ファクトチェック・イニシアティブ(以下 FII) が提供する FactCheck Navi を使用する。2021 年 4 月現在で 600 超件のファクトチェック結果が公表されている。一方ファクトチェックにより正確と判断された事例はフェイクに比べて少ないため、正確なニュースとして大手新聞社やロイター通信等の記事を収集する。

また目標Ⅱに向けて正解ラベルがなく、弱いアノテーションを付加する対象である記事を追加する。正確とみられるニュースは先述と同じく大手マスメディアが発信したニュースを扱い、疑わしい記事として FactCheck Navi によって虚偽と 3 回以上判断されたことがあるニュースサイトの他記事も収集対象とする。最終的には真偽合わせてラベル付き記事を約 1200 件、ラベルなし記事を約 5000 件以上収集を目指す。ユーザの反応として、全記事を対象に SNS 上で寄せられたコメントとして Twitter にて記事 URL を含むツイートも収集する。

**目標Ⅱ: 弱教師あり学習によってラベル不足を補うモデルを構築する (1 年目 - 2 年目)**

ラベル付きユーザの反応を弱いアノテーションとして扱ってフェイクニュースを検出する方法は、英文記事を対象にした Kai Shu らの研究で 1

目標	採用前	1年目	2年目
日本語データセットの作成	←	→	
自動検出モデルの構築		←	→
ユーザへの提供システム開発			←

図3: 研究計画の年毎の流れ。

例が示されている [9]。ここではコメント群の感情値の標準偏差やコメント者の過去の投稿、そしてフォ

(研究目的・内容等の続き)

ロー関係から弱いアノテーションを付加している。これら3種の弱いアノテーションも併せて学習することで、推論時はユーザの反応を使わずに正確な早期検出を実現した。日本語で実現を目指す場合、日本語での検出を行う研究が少なく言語の違いによる影響が未知数であるため大幅なモデルの改変を要する可能性がある。今回は日本語での早期検出実現に向けて3種類の弱いアノテーションに加え、投稿者のプロフィールや使用された絵文字、ハッシュタグといったコメント情報で有用なものがないか模索する。実験では、学習時には記事とユーザの反応を入力に扱い、テスト時は記事のみを入力して早期検出時の状況を再現し既知の手法と比較する。既知の手法と比較して検出性能の改善がみられた場合は成功とみなす。

### 目標III: 利用者とのやり取りで検出を行うモデルを開発する (2年目)

データセットとモデルの構築のみでは、今後のニュース内容の変化に対応することが難しい。継続して社会の潮流の変化に対応するためには、データセットと検出方法を工夫する必要がある。また、本研究の最終的な目標は国連の Pause/ちょっと待って運動 [10] と同じく SNS ユーザに対してフェイクニュース拡散前にユーザに再考を促して抑止を行うことである。この実現には、検出モデルの構築のみならずユーザに対し納得できる形で記事の疑わしさを提供するシステムの構築により初めて実現する。本研究は、ユーザがモデルに直接真偽を問い合わせるシステムを構築する。上記で実装したモデルに対して、ユーザが記事内容を入力し問い合わせることで、真偽判断結果を提供すると共にデータセットへラベルなしニュースの追加を狙う。実験では SNS 利用者を対象とした主観評価実験により、フェイクニュース拡散の抑止効果を測定する。モデル提供が難しい場合は、データセットとモデルを公開することで日本語におけるフェイクニュース早期検出研究の促進を目指す。

### ③研究の特色・独創的な点

#### 本研究の特色

- 日本語を対象にフェイクニュースの自動検出を行う点。
- ファクトチェックの結果を待たず早期の検出を目指す点。
- ユーザからの情報提供で継続してデータセット拡大・モデル改善を行う点。

#### 独創的な点: 先行研究との比較

深層学習でフェイクニュースを自動検出する研究対象は英文に集中しており、日本語データセットがない。言語に囚われずユーザによる拡散された経緯で真偽を判断する研究もあるが [11]、依然として記事の内容を考慮した研究では日本語を対象としたものがない。

#### 独創的な点: 予想されるインパクト・将来の見通し

総務省によると SNS 利用率は 2019 年現在 69% を占める上、SNS マーケティング市場規模は 2025 年に 1 兆 1,171 億円まで成長する (出典: サイバー・バズ/デジタルインファクト調べ) と推測されている。SNS 利用者が拡大を続ける中、本研究の完成によりこれまで活発になされていなかった日本語のフェイクニュースを早期検出するモデルの開発および提供が可能となる。SNS 利用者への注意喚起に活用ができるほか、ファクトチェックを行う人への補助システムへの活用といった様々な形式で SNS 上で騙される人が増え社会的損害や風評被害の発生を未然に防ぐ仕組み作りに貢献する可能性がある。

### ④申請者が担当する部分

本研究は所属研究室内でも萌芽的な取り組みで、環境・技術面の支援を除き申請者が全部分を担当する。データセットの生成では、正確なニュースの取得へ大手マスメディアへ協力を求める可能性がある。

### ⑤受入研究機関と異なる研究機関での研究従事計画

申請者は1年間タリン工科大学の言語技術研究所 (Tanel Alumäe 所長) で活動予定である。当該分野は北米と欧州の研究が活発であることから、最前線の研究に従事するために必要である。

参考文献

[9] K. Shu, et al. ECML-PKDD 2020

[11] Tarek H., et al. ICDCIT 2020

[10] UNIC. [https://is.gd/UNIC\\_pause](https://is.gd/UNIC_pause). 2020

### 3. 人権の保護及び法令等の遵守への対応 ※本項目は1頁に収めてください。様式の変更・追加は不可。

本欄には、「2. 研究計画」を遂行するにあたって、相手方の同意・協力を必要とする研究、個人情報の取り扱いの配慮を必要とする研究、生命倫理・安全対策に対する取組を必要とする研究など法令等に基づく手続が必要な研究が含まれている場合に、どのような対策と措置を講じるのか記入してください。例えば、個人情報を伴うアンケート調査・インタビュー調査、国内外の文化遺産の調査等、提供を受けた試料の使用、侵襲性を伴う研究、ヒト遺伝子解析研究、遺伝子組換え実験、動物実験など、研究機関内外の情報委員会や倫理委員会等における承認手続が必要となる調査・研究・実験などが対象となりますので手続の状況も具体的に記入してください。

なお、該当しない場合には、その旨記入してください。

コメント取得を予定してしている SNS は Twitter である。Twitter 社は 2020 年 3 月より学術目的で Twitter API の利用を自由化しているほか、取得したツイート ID を含む情報をデータセットとして公開することも学術目的であれば認められている [12]。

また、先行研究が提供したデータセットを使用する場合は、提供者が示しているライセンスやポリシーを遵守する。

なお、学習済みモデルの公表は平成 30 年改正著作権法第 30 条 4 号により認められている。

ただし、本研究では主観評価実験として SNS ユーザを対象としたアンケート調査を予定している。この調査により収集したデータは、個人の特定につながる情報を匿名化した上で解析を行い、解析結果の公表に際しては、匿名化を行ったデータを用い、個人情報の漏洩防止に配慮する。

#### 参考文献

- [12] Twitter 開発者ポリシーを分かりやすくアップデート, 2020 年 3 月 11 日. (最終閲覧日 2020 年 4 月 19 日) [https://blog.twitter.com/developer/ja\\_jp/topics/tools/2020/DevPolicyUpdate.html](https://blog.twitter.com/developer/ja_jp/topics/tools/2020/DevPolicyUpdate.html)

**4. 【研究遂行力の自己分析】** ※各事項の字数制限はありませんが、全体で2頁に収めてください。様式の変更・追加は不可。

本申請書記載の研究計画を含め、当該分野における(1)「研究に関する自身の強み」及び(2)「今後研究者として更なる発展のため必要と考えている要素」のそれぞれについて、これまで携わった研究活動における経験などを踏まえ、具体的に記入してください。

**(1) 研究に関する自身の強み**

自ら抱いた問題意識を出発点に研究を行う主体性

貪欲な海外論文調査に裏打ちされた状況把握能力

産学問わない活動で培ったプログラミング能力で実現される実装能力

プラットフォームを問わず議論を活発に行えるコミュニケーション能力

相手が小学生でも物事を分かりやすく伝えることができるプレゼン能力

**(2) 今後研究者として更なる発展のため必要と考えている要素**

要素 1: 学術的成果と社会問題の最前線の間にあるギャップを埋めるための発想と問題解決力

要素 2: 多彩な分野や言語・地域圏の研究者らと活発な議論を交わす能力

要素 3: 研究で得られた成果をどんな聞き手でも分かりやすく伝えられる表現力

(研究遂行力の自己分析の続き)

**5. 【目指す研究者像等】** ※各事項の字数制限はありませんが、全体で1頁に収めてください。様式の変更・追加は不可

日本学術振興会特別研究員制度は、我が国の学術研究の将来を担う創造性に富んだ研究者の養成・確保に資することを目的としています。この目的に鑑み、(1)「目指す研究者像」、(2)「目指す研究者像に向けて特別研究員の採用期間中に行う研究活動の位置づけ」を記入してください。

**(1) 目指す研究者像** ※目指す研究者像に向けて身に付けるべき資質も含め記入してください。

自分の興味のある分野を研ぎ究めると同時に社会問題を解決して人々の生活を幸せにしたい

嘘の情報に騙されて誤った風評が残り不幸になる人を0にしたい

ファクトチェックでは嘘は嘘であると騙された人を相手に分かりやすく説明することが重要である

自己完結のみならず成果を他人に伝えるまでが研究である

**(2) 上記の「目指す研究者像」に向けて、特別研究員の採用期間中に行う研究活動の位置づけ**

特別研究員の採用期間中に行う研究活動のなか、4-(2)で挙げた今後研究者としてさらなる発展のため必要と考えている要素の習得を通して、学術研究で得た知見を直接 SNS ユーザを含む日本社会に還元する研究者を目指す。その実現に向け、査読付き国際会議ないしは国際論文誌への論文発表をはじめ、国内・国際会議での口頭発表も積極的に行う。また、自然言語処理コミュニティに限らず国内ニュースメディアと積極的に連携を行い、フェイクニュースの自動検出に関連した共同研究の実現が理想である。

特別研究員として研鑽を重ねていき、現状の研究への新たな発想を追加し、実現に向けて幅広い人々と議論を重ね、得られた成果を端的に説明することが、能動的に一貫して社会課題を解決へ自ら導く研究者として大成する。その実現の大きな足がかりが本研究計画である。