



プロジェクトの更新があればここに追記予定

ERUSTA

Extension that Recommends Useful and Secure Technical Articles

学習駆動コース坂井ゼミ

竹内 悠人

～ 有用な技術記事を推薦する拡張機能～

ERUSTAとは？

～システム構成～

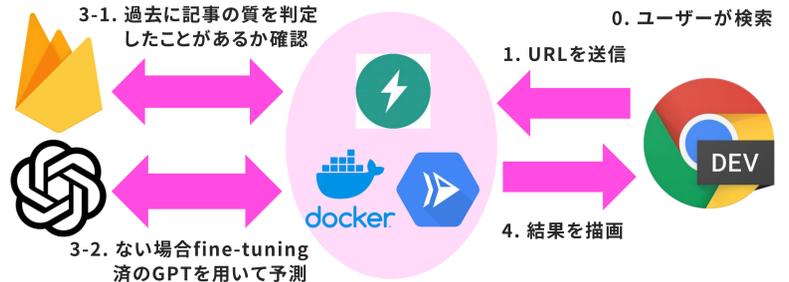
有用性の低い記事たち

- ・アフィリエイト目的
- ・タイトル詐欺
- ・内容の薄い記事
- ・危険なドメイン

削除

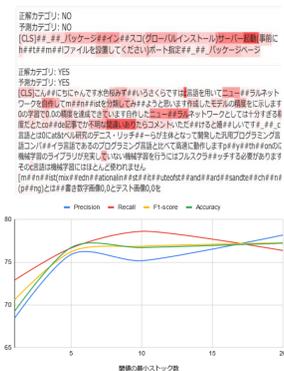


- ・質の高い記事のみを表示
- ・(将来的には) 個々や団体の技術力に合わせた推薦



～SecHack365での一年～

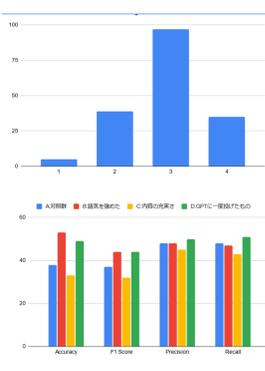
実験1: 記事の良し悪しって機械にわかるの？



- ・初期のデータセットとして、Qiitaのストック(お気に入り)数を利用
- ・2万件の記事を2種のラベルに分類して、データセットを作成
- ・このデータセットで日本語BERTモデルを学習しAttentionを可視化
- ・最終スコアは、最小ストック数の閾値を20にした時が一番高かった(77.3%)
- ・Attentionの可視化により、記事の有用性を特徴づける文章自体の抽出は難しいことがわかった

⇒精度に限界! 有用さはストック数だけでは判断できないが、記事の良し悪しはいい感じに学習できそう

実験2: 良し悪しの判断を人間に近づけられるか



- ・実験にエンジニア15人が参加
- ・ランダムに生成した20個の検索ワードを用意
- ・検索ワードでの結果上位10件の記事を収集
- ・各記事ごとに有用度を3人のエンジニアが2段階で評価
- ・左の4つの条件でGPT-3.5 Turboを学習
- ・結果どれも性能が伸び悩んだ(50%)

⇒二値分類に帰着したのに精度がおしまい 厳しい

実験3: 拡張機能のプロトタイプの利用者ビリティ調査



- ・拡張機能のプロトタイプを作成した
- ・実験に15人が参加
- ・ユーザービリティをSUSを用いて評価した
- ・システムのシンプルさや操作性で高評価を得た
- ・Q1, Q5, Q6, Q9が低かった一番の原因は、品質の推論結果が表示されるまでに時間がかかってしまったからと考えられる

⇒推論に時間がかかってしまうことが問題

反省① 推論に時間がかかってしまっている

⇒推論を並列化 + 過去に一度推論した記事はDBに保存

反省② 作成したデータセットのサイズが小さい

⇒検索キーワードを大量に用意(客観性保持) + 評価軸を増やす

“有用性”を再考する

- ・有用性を特徴づける要素は多種多様であり、どれを評価軸として用いるのが適切なかわからない
- ⇒先行研究をあたりたい
- ・先行研究:記事の有用性の特性をおおまかに、「戦略的」「機能的」「体験的」に分けた時、記事の有用性を特徴づけるのに良く使われるパラメータは120個存在する(右図)[1]
- ⇒この中から技術記事での有用性に使えるものを選定
- ・加えて技術記事に特化させる目的でパラメータを追加した

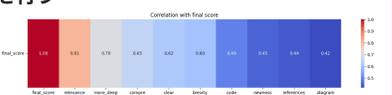
| Dimension | Strategic | Functional | Experiential |
|-----------------------------|---|--|---|
| Usability and accessibility | Effectiveness Efficiency Customization Time prevention | Availability Breadth of content Space saving Visibility | User control Learnability Ease of comprehension Flexibility Interactivity Recognition labels |
| Content and services | Pragmatic discourse Equity and inclusion Expertise Content information Social responsibility Social responsibility Tutorials and expert | Up-to-date Attribution of authorship Written for the web Writing and editing Lack of content Variation of information Timeliness and rigor | Technical character Help and support Completeness and robustness Reliability Flexibility Personalization |
| Information architecture | Discovery or predictive recommendations Home page or main page | Internal site search Preventing repeat steps Clear navigation Short steps Consistency Utility | Internal site search Mobile usability Localization and state Linking Personalization Versatility |
| User experience | Trust Satisfaction Perceived value | Consistency Credibility Expertise Empathy | Consistency Credibility Expertise Empathy |

- 機能的
 - コンテンツとサービス
 - 最新の情報が反映されているか?
 - 参考・引用文献が書かれているか?
 - 図表が用いられているか?
- 体験的
 - ユーザビリティとアクセシビリティ
 - より深い学びに繋がったか?
 - 明確に記事の内容を理解できたか?
 - 記事は簡潔であったか?
 - コンテンツとサービス
 - 記事の中身は網羅性があったか?

検証

- ・情報系の大学生4人+エンジニア17人がこの1200件の記事を先ほど紹介した9個の有用性を特徴づける要素について5段階での点数付けを行った
- ・検索キーワード、記事の本文、検索結果の順位のみを用いてGPT3.5のFine-tuningを行うと3段階、5段階評価でそれぞれF1scoreはそれぞれ38%, 19%
- ⇒そのままでは有用な記事の推薦は困難
- ・有用性を特徴づける9個の要素を追加して同様の操作を行うと、F1scoreはそれぞれ74%, 52%
- ⇒高確率で良質な記事を推薦することが可能!
- ・また最終的なスコアとの相関が高い要素を調べたところ、関連性、より深い理解につながるか、網羅性があったかなどが高い相関をもっていることが確かめられた

| 項目 | 検索 | 本文 | 順位 |
|----|-----|-----|-----|
| Q1 | 0.8 | 0.7 | 0.6 |
| Q2 | 0.7 | 0.6 | 0.5 |
| Q3 | 0.6 | 0.5 | 0.4 |
| Q4 | 0.5 | 0.4 | 0.3 |
| Q5 | 0.4 | 0.3 | 0.2 |
| Q6 | 0.3 | 0.2 | 0.1 |
| Q7 | 0.2 | 0.1 | 0.0 |
| Q8 | 0.1 | 0.0 | 0.0 |
| Q9 | 0.0 | 0.0 | 0.0 |



クエリ作り直し

- ・タイトルからクエリを作りやすい
- ・質問にいいねを付ける文化があるため質問の質を評価できることから日本語Stack Overflowを採用
- ・右に示す人気タグ40個から上位5件の質問のタイトルを検索クエリ化
- ・先ほど作成した検索クエリをGoogle検索に投げ、上位3件+18,19,20番目の記事を取得(質の悪い記事も集めるため)
- ⇒40テーマ×5クエリ×6記事 = 1200記事を取得

| id | topic | query |
|----|---------------|----------------------|
| 1 | python | 5 manaca |
| 2 | python | 5 manaca |
| 3 | javascript | 5 manaca |
| 4 | java | 5 onjective-c |
| 5 | php | 5 pandas |
| 6 | swift | 5 aws |
| 7 | ruby_on_rails | 5 aws |
| 8 | html | 5 unity |
| 9 | c# | 5 sql |
| 10 | ios | 5 git |
| 11 | android | 5 docker |
| 12 | ruby | 5 modis.jp |
| 13 | c++ | 5 centos |
| 14 | linux | 5 visual studio |
| 15 | kcode | 5 laravel |
| 16 | mysql | 5 algorithm |
| 17 | windows | 5 st |
| 18 | css | 5 google-apps-script |
| 19 | c | 5 json |
| 20 | jquery | 5 .NET |
| 21 | python3 | 5 regular expression |

引用・謝辞

[1]Alejandro Morales-Vargas-2023
Website quality evaluation model for developing comprehensive assessment instruments based on key quality factors
Table 6. Most frequently mentioned website quality parameters, organized by focus

実験1は国立研究開発法人科学技術振興機構グローバルサイエンスキャンパス(GSC)「情報科学の達人」育成官員協会のプログラムの支援のもと実施しました。また実験2を支援していただいた株式会社ラックサイバーセキュリティラボ「ITスペシャリスト」やラボ「ITスペシャリスト」の方々にも多大なる感謝を申し上げます。

外部での発表(IPSJ2024)

RAGを用いた有用性の高い技術記事を推薦するモデルとデータセットの開発
竹内悠人・菅原期