

ニューラルネットワークによるTwitterエンゲージメントの予測

石川県立金沢泉丘高等学校

研究の流れ

- ① TwitterAPIを用いてツイートを集める
- ② ツイートを絵文字の消去などの前処理にかける
- ③ 前処理済みのツイートを分かち書きする
- ④ Doc2Vecを用いてツイートをベクトル(数値データ)に変換する
- ⑤ ベクトルとフォロワー数・フォロー数・曜日・時刻のデータを使ってモデルを作成し、検証・考察する

モデル作成の下準備

利用システム

TwitterAPI

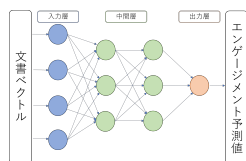
Twitter社から提供されている特定のワードを含んだツイートを収集することができるもの

Doc2Vec

ライブラリ「Gensim」で使用できる文章をベクトルに変換することができる機能[1]

Neural Network

入力した値に、重みをかけて線形変換した値を次の層へ送る、という計算操作がネットワーク状に連なったもの。ネットワーク全体の重みを調整する操作1回を1epochと呼ぶ。調整方法はAdamというものを用いた。



イメージ図

評価関数

平均二乗誤差(MSE)

モデルの学習時にMSEの値が小さくなることを目的とする。

$$MSE = \frac{1}{n} \sum_{k=i}^n (y_i - \hat{y}_i)^2$$

y_i : i 番目のデータの真の値
 \hat{y}_i : i 番目のデータの予測値

予測対象とするエンゲージメントはいいね数とリツイート数の和として計算しておく

パターン

ニューラルネットワークに入力するデータの組み合わせ

パターン①: テキスト

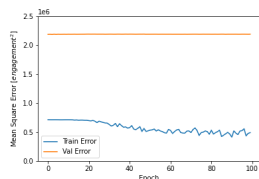
パターン②: テキスト・フォロワー数・フォロー数・曜日・時刻

パターン③: フォロワー数・フォロー数・曜日・時刻

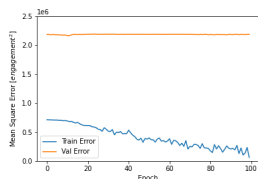
モデルの作成と考察

「オススメ」という言葉の含んだツイートを22029件収集し、モデルを学習した。ベクトルは400次元とした。

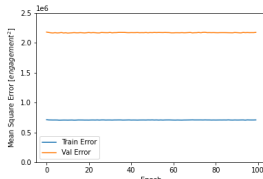
パターン①



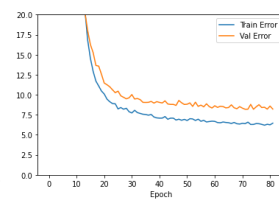
パターン②



パターン③



- ・縦軸: MSE 横軸: epoch
- ・オレンジ: 検証データ
- ・青色: 訓練データ
- ・正しく学習できているときのグラフは右上のような推移になる[2]



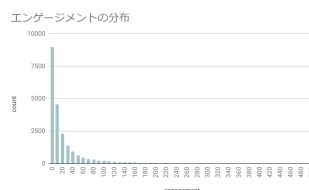
- ・学習しても未知のデータである検証データに対して予測の誤差が減少していない(収集データの調査へ)
- ・ベクトルの次元数が大きすぎて入力が多くなり予測がうまくいっていない可能性がある(モデルの改善へ)
- ・パターン①とパターン②より、テキスト以外のデータが予測に影響を及ぼしていると考えられる
- ・パターン③では訓練データでも学習が進まないためテキストは確実に予測に必要なとわかる
- ➡ただし、一般的な感覚ではテキスト以外のデータとエンゲージメントの間だけでも相関性がありそう(モデルの改善へ)

収集データの調査

今回使用した「オススメ」という言葉を含んだツイート22029件のエンゲージメントのヒストグラムを作成した。



エンゲージメント250以下のものが大多数



22029ツイート中エンゲージメントが500以上のものは313ツイートのみ

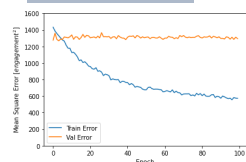
外れ値と捉えられる

エンゲージメントが極端に大きいツイートが大きな外れ値として予測に影響を及ぼしている可能性がある

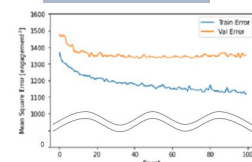
モデルの改善

エンゲージメントの範囲を250以下に、ベクトルの次元を100次元にして再度パターン②、③で学習を行った。

パターン②



パターン③



- ・パターン②は正しく予測できていない
- ➡ベクトルの次元がまだ大きすぎるかベクトル変換がよくない
- ・パターン③はグラフが下がった
- ➡エンゲージメントが250以下の小さいツイートではフォロワー数・フォロー数・曜日・時刻とエンゲージメントには少なくとも相関がある

今後の展望

- ・より高精度なモデルの作成
- ・どのような言葉の含んだツイートでも予測することのできる一般的なモデルの作成
- ・多くのエンゲージメントを得る方法の理論化と実社会で広告への応用

参考文献

- [1]Tomas Mikolov. Distributed Representations of Sentences and Documents. https://cs.stanford.edu/~quocle/paragraph_vector.pdf(参照2021-11-10).
- [2]Google. "帰帰: 燃費を予測する" <https://www.tensorflow.org/tutorials/keras/regression?hl=ja>(参照2021-11-10)他多数

抄録

総務省の調査では、2021年1月時点で日本人口の42.3%がTwitterを利用しているため、ツイートに対して多くの人々から反応を得られる可能性がある。そこでニューラルネットワークを用いて反応を表すエンゲージメント（いいね・リツイート数）を予測するモデルを作成することにした。

1. 研究の背景と目的

多くの人々がTwitterを利用しているため、ツイートに対して様々な反応が存在する。そこでツイートがどれだけのエンゲージメントを得られるかを予測するモデルの作成を目指す。

2. 方法

1. TwitterAPIを用いてツイートデータを収集する。
2. テキスト内の重ね表現や、数、絵文字に関して処理を行う。
3. Doc2Vecを用いて、テキストを数値化(ベクトル化)する。
4. そのベクトルにツイートの情報を付加したデータを、予測の基とする訓練用データと予測の誤差を測るテスト用データに分割し、前者を用いていいね・リツイート数を予測するニューラルネットワークモデルを作る。
5. テスト用データを用いて予測の誤差を測定する。
6. ツイートの情報であるテキスト・フォロワー・フォロワー・時刻の組み合わせを変更し、より予測の誤差が小さいモデルを作成する。
- 7.

3. 結果

様々なパターンを試したが、すべてにおいていいね・リツイート数の誤差の平均がおおよそ80となった。また、訓練用データのみ誤差の減少が見られた。

4. 考察

訓練用データのみで特化したモデルができてしまったことや、付加する情報を変化させても誤差の変化がほぼなかったことより、モデルの訓練の方法または使用するデータに問題があると考えられる。

5. 結論

作成したモデルは誤差おおよそ80であり、訓練方法や使用データによってより誤差の小さいモデルの作成ができる可能性がある。

6. 参考文献

Tomas Mikolov. Distributed Representations of Sentences and Documents.

https://cs.stanford.edu/~quocle/paragraph_vector.pdf (参照2021-11-10).

Google. ” 回帰：燃費を予測する”

<https://www.tensorflow.org/tutorials/keras/regression?hl=ja> (参照2021-11-10).

7. キーワード

TwitterAPI Doc2Vec ニューラルネットワーク 自然言語処理