

進捗報告

1 今週やったこと

- BERT-Transformer モデルで多値分類をした.
- BERT-Transformer モデルに, 各ラベルと関連度の高いフレーズを抽出するターゲットフレーズラベリング器の実装中

2 データセット

2.1 楽天トラベルレビュー : アスペクトセンチメントタグ付きコーパス

楽天トラベルレビュー : アスペクトセンチメントタグ付きコーパスとは, 楽天グループ株式会社が提供しているデータセットである. 日本語レビュー文章とそれぞれの文章について, 立地, 部屋, 食事等の 7 項目のアスペクトに対するポジティブまたはネガティブのタグが付与されている. 総データ数は 76624 で, そのうち全ラベルが 0 であるデータは 23432 である. 今回は 14 のいずれのラベルにも属さないデータを除くことで, 少なくとも 1 つのラベルに属し, 語彙数が 10 以下と 100 以上のデータを取り除いたデータ群にした. 総データ数は 50211 である. 表 1 にデータの具体例を示す.

3 実験

3.1 実験 1

楽天トラベルレビュー : アスペクトセンチメントタグ付きコーパスを用いて多値分類をした. 訓練データ数を 11200, バリデーションデータを 2800 としして学習をした. データに含まれるラベルが 0 が多く, 評価指標として正解率を用いると高くなってしまい正しく評価できないと考えたので F1 値を用いた. 表 2 に学習時のパラメータを示す. 2 値分類の時と異なるのはデータセットや損失関数である.

表 2: BERT-Transformer モデルの実験時のパラメータ

パラメータ	値
BERT 層の入力次元数	768
BERT 層の出力次元数	768
Transformer 層の層数	1
Transformer 層の入力次元数	768
Transformer 層の出力次元数	14
バッチサイズ	14
最適化関数	Adam
学習率	0.0001
損失関数	BCEWithLogitsLoss
エポック数	30

続いて図 1 に訓練時の訓練データとバリデーションデータの損失を示す.

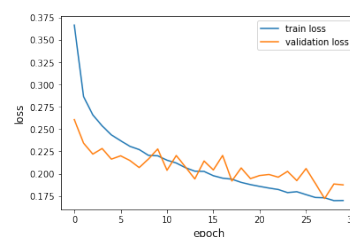


図 1: 実験時の訓練データとバリデーションデータの正解率の推移

また, 図 2 に実験 1 における学習時のバリデーションデータの F1 値の推移を示す.

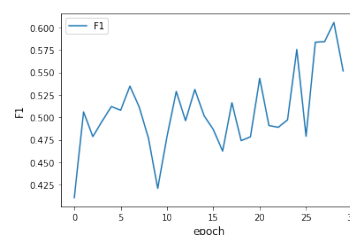


図 2: 実験時の F 値の推移

図 1 におけるバリデーションデータの損失の推移や図 2 における F 値の推移に加えて, 精度に問題を感じたため, 改善を図るためにプログラムの修正をしています.

表 1: 両方のラベルが立っているデータの具体例

テキスト	朝食 po	朝食 ne	夕食 po	夕食 ne	風呂 po	風呂 ne	サービス po	サービス ne	立地 po	立地 ne	設備 po	設備 ne	部屋 po	部屋 ne
立地、最上階、部屋からの景色、エアウィーヴ等の 良い点と比較しても、次の機会に泊まるかは疑問です.	0	0	0	0	0	0	0	1	1	0	1	1	1	1
外観を見て失敗したと思いましたが、中に入ると 別世界でした.	0	0	0	0	0	0	0	0	0	0	1	1	0	0
古いながらも大変メンテナンスされていますので 清潔でした.	0	0	0	0	0	0	0	0	0	0	1	1	1	1
食事も夕・朝とも質量ともに問題なかったのですが、 逆に朝は量が多すぎるくらいでした.	1	1	1	0	0	0	0	0	0	0	0	0	0	0