

データサイエンス特論 プログラム課題 第2回(多値ラベル分類)

以下の課題を実施し、実行したプログラムと実行結果(プログラムはPython, Jupyter notebook(拡張子.ipynb, .html)等でログをまとめるのも可、実行結果をMS Word(PDF)等でまとめることも可)を、ZIPでまとめ、MoodleLMSにアップロードすること。期限は【8月8日(月)】の深夜23:59分までとする。(プログラム課題はこれで終わりです)

UCI(カリフォルニア大学アーバイン校)の機械学習用データレポジトリにある有名なロイターのニュース記事(Reuters-21578) <https://archive.ics.uci.edu/ml/datasets/reuters-21578+text+categorization+collection> は、英語で書かれた1980年当時のニュース記事の集合である。このデータは、機械学習の練習にしばしば利用されてきた。ニュース記事にはカテゴリ情報が1種類以上付与(マルチラベルが付与)されており、このカテゴリラベルを予測するテキスト分類問題を考える。ただし、元のロイターニュース記事には、各種のノイズ文書がある。そこで、55カテゴリ、10,700文書(7713訓練文書, 2987テスト文書)に絞り込んだ部分集合記事データを作成した。以下のサイトからダウンロードできる。 https://www.kde.cs.tut.ac.jp/~aono/data/ma_reuters.zip においてある。以降、このデータをMA-Reutersと呼ぶこととする。

作業手順を述べる。

(Step 1) MA-Reutersを各自の作業できるPC等にダウンロードせよ。

(Step 2) MA-ReutersはPythonのNLTKパッケージのコーパスに従って作成してある。そこで、各自のPythonの環境にNLTKがまだインストールされていない場合、NLTKをインストールせよ。同時に <https://www.nltk.org/data.html> を参照してデータのダウンロードを最初に一回行い、適当なNLTK用のデータフォルダーを作成しておくこと。たとえばWindows環境下で、C:\nltk_dataに作成したとすると、コーパスフォルダが、C:\nltk_data\corpora以下にできるので、(1)で入手したma_reuters.zipをcorporaフォルダにコピーしておくこと。

(Step 3) NLTKと同時にscikit-learn(サイキットラーン)パッケージも利用することになるので、こちらのインストールがまだの場合しておくこと。

以下が課題である。本課題で Python のパッケージのインストールやニューラルネットワーク (NN) の環境設定などのヘルプはできないことを注記する。

- [1] <https://www.kde.cs.tut.ac.jp/~aono/2022/DataScience/Reuters-Multi-Label-SVM-Example.html> にある Jupyter notebook での Python コードを順番に実行し、多値クラス多値ラベル問題を解決できることを確認せよ。ただし、このコード内でオレンジカテゴリー内の(最初の)文書のプリントは、違うカテゴリー (オレンジカテゴリー以外) の文書に変更して実行すること。
- [2] 55 個のカテゴリーのなかで、Jaccard 係数 (資料第 6 回参照) が 最も高いカテゴリーと最も低いカテゴリーが何であったかのべよ。
- [3] ここでの NLTK の TF-IDF モデル+SVM 以外の組み合わせで本マルチラベル問題に対応できる手法を適宜ためし、実行結果をのべよ。たとえば、テキスト表現に授業の資料やビデオで述べた Word2Vec や BERT 等を用いてもよい。NLTK の BoW(TF-IDF)モデルのまま場合、少なくとも分類器は変更すること。たとえば、ニューラルネットワーク(NN)を用いてもよい。NN の場合、評価手法はバイナリークロスエントロピー等の誤差がマルチラベル・マルチクラス分類ではよく用いられる。