

第6章: 機械学習

本章では、Fabio Gasparetti氏が公開している[News Aggregator Data Set](https://archive.ics.uci.edu/ml/datasets/News+Aggregator) (<https://archive.ics.uci.edu/ml/datasets/News+Aggregator>)を用い、ニュース記事の見出しを「ビジネス」「科学技術」「エンターテインメント」「健康」のカテゴリに分類するタスク（カテゴリ分類）に取り組む。

50. データの入手・整形

[News Aggregator Data Set](https://archive.ics.uci.edu/ml/datasets/News+Aggregator) (<https://archive.ics.uci.edu/ml/datasets/News+Aggregator>)をダウンロードし、以下の要領で学習データ（`train.txt`）、検証データ（`valid.txt`）、評価データ（`test.txt`）を作成せよ。

1. ダウンロードしたzipファイルを解凍し、`readme.txt`の説明を読む。
2. 情報源（`publisher`）が“Reuters”, “Huffington Post”, “Businessweek”, “Contactmusic.com”, “Daily Mail”の事例（記事）のみを抽出する。
3. 抽出された事例をランダムに並び替える。
4. 抽出された事例の80%を学習データ、残りの10%ずつを検証データと評価データに分割し、それぞれ `train.txt`, `valid.txt`, `test.txt` というファイル名で保存する。ファイルには、1行に1事例を書き出すこととし、カテゴリ名と記事見出しのタブ区切り形式とせよ（このファイルは後に問題70で再利用する）。

学習データと評価データを作成したら、各カテゴリの事例数を確認せよ。

51. 特徴量抽出

学習データ、検証データ、評価データから特徴量を抽出し、それぞれ `train.feature.txt`, `valid.feature.txt`, `test.feature.txt` というファイル名で保存せよ。なお、カテゴリ分類に有用そうな特徴量は各自で自由に設計せよ。記事の見出しを単語列に変換したものが最低限のベースラインとなるであろう。

52. 学習

51で構築した学習データを用いて、ロジスティック回帰モデルを学習せよ。

53. 予測

52で学習したロジスティック回帰モデルを用い、与えられた記事見出しからカテゴリとその予測確率を計算するプログラムを実装せよ。

54. 正解率の計測

52で学習したロジスティック回帰モデルの正解率を、学習データおよび評価データ上で計測せよ。

55. 混同行列の作成

52で学習したロジスティック回帰モデルの混同行列（confusion matrix）を、学習データおよび評価データ上で作成せよ。

56. 適合率，再現率，F1スコアの計測

52で学習したロジスティック回帰モデルの適合率，再現率，F1スコアを，評価データ上で計測せよ。カテゴリごとに適合率，再現率，F1スコアを求め，カテゴリごとの性能をマイクロ平均（micro-average）とマクロ平均（macro-average）で統合せよ。

57. 特徴量の重みの確認

52で学習したロジスティック回帰モデルの中で，重みの高い特徴量トップ10と，重みの低い特徴量トップ10を確認せよ。

58. 正則化パラメータの変更

ロジスティック回帰モデルを学習するとき，正則化パラメータを調整することで，学習時の過学習（overfitting）の度合いを制御できる。異なる正則化パラメータでロジスティック回帰モデルを学習し，学習データ，検証データ，および評価データ上の正解率を求めよ。実験の結果は，正則化パラメータを横軸，正解率を縦軸としたグラフにまとめよ。

59. ハイパーパラメータの探索

学習アルゴリズムや学習パラメータを変えながら，カテゴリ分類モデルを学習せよ。検証データ上の正解率が最も高くなる学習アルゴリズム・パラメータを求めよ。また，その学習アルゴリズム・パラメータを用いたときの評価データ上の正解率を求めよ。