# Unsupervised Text Classification

Iqbal Pahlevi Amin \*
Faculty of Computer Science
Universitas Indonesia
Depok, Indonesia
Email: \*iqbal.pahlevi@ui.ac.id

Abstract—Pada eksperimen ini, penulis melakukan klasifikasi dokumen tanpa label berdasarkan label yang telah didefinisikan di awal melalui pendekatan unsupervised learning. Pendekatan ini memerlukan label beserta deskripsi yang telah didefinisikan di awal serta dataset tanpa label yang akan digunakan. Metode yang digunakan bekerja dengan cara melakukan embedding vektor pada dokumen dan deskripsi masing-masing label. Kemudian, masing-masing vektor embedding dokumen akan dihitung similarity-nya dengan vektor embedding semua deskripsi label. Penulis mendapatkan F1-score terbaik sebesar 0.823 pada data training dan 0.817 pada data testing.

Index Terms—Unsupervised learning, vektor embedding

#### I. PENDAHULUAN

Paper ini membahas tentang cara melakukan klasifikasi teks atau dokumen tanpa label melalui pendekatan *unsu-pervised* learning dengan menggunakan Lbl2Vec. Klasifikasi teks atau dokumen umumnya dilakukan melalui pendekatan konvensional, yaitu *supervised learning*. Namun, pendekatan ini membutuhkan jumlah *labeled data* yang sangat besar untuk *training* model. Padahal pada kenyataannya, *labeled data* dalam bentuk dokumen itu sering tidak tersedia. Di sisi lain, apabila ingin melakukan labeling data secara manual, itu akan memakan banyak waktu dan membutuhkan *resource* yang banyak. Oleh karena itu, pendekatan *unsupervised learning* mulai digunakan karena lebih efisien dan murah. Klasifikasi teks secara *unsupervised* disebut juga sebagai *zero-shot text classification*.

pada eksperimen ini penulis mereproduksi eksperimen pada paper[1]. Tidak seperti pada paper acuan, pada eksperimen ini penulis melakukan eksperimen dengan hanya menggunakan dataset dari AG's corpus tanpa data dari 20Newsgroups. Penulis melakukan menjalankan tiga skenario dalam melakukan training model. Skenario ini berbeda dalam jumlah data training yang akan digunakan. Skenario pertama, menggunakan 10% data training. Skenario kedua menggunakan 25% data training. Dan skenario ketiga menggunakan 100% data training. Dari ketiga skenario tersebut didapatkan bahwa semakin banyak data training yang digunakan, model dapat menghasilkan nilai F1-score yang lebih bagus. Pada penggunaan 100% data training, penulis mendapatkan F1score yang tidak jauh berbeda dari hasil eksperimen pada paper acuan[1]. Selain itu, penulis juga melakukan eksperimen dengan melakukan empat percobaan dengan inisiasi model menggunakan hyperparameter yang berbeda.

Selanjutnya pada bab alur kerja program akan dijelaskan lebih detail terkait alur program berjalan, dari proses tokenisasi

hingga proses evaluasi pada data testing. Kemudian, pada bab percobaan dan analisis hasil akan dijelaskan terkait setup eksperimen dan penjelasan terkait hasil yang didapatkan. Terakhir, pada bab kesimpulan dan saran akan dijelaskan kesimpulan yang didapat penulis setelah melakukan eksperimen ini serta beberapa saran yang dapat dilakukan pada penelitian selanjutnya.

#### II. CARA KERJA PROGRAM

Secara garis besar, program ini berjalan melalui empat tahap penting. Pertama, melakukan *data preprocessing*. Kedua, melakukan *embedding* vektor untuk label dan dokumen *training-testing*. Ketiga, mengategorikan masing-masing dokumen *training* ke suatu label berdasarkan hasil perhitungan *similarity*. Dan terakhir, mengategorikan masing-masing dokumen *testing* ke suatu label berdasarkan hasil perhitungan *similarity*.



Fig. 1. Dokumen train dan test beserta hasil tokenisasinya.

Gambar 1 menampilkan cuplikan data *training* dan *testing* beserta hasil tokenisasinya. Proses tokenisasi data memanfaatkan fungsi simple\_preprocess yang berasal dari *library* gensim. Hasil tokenisasi kedua jenis data disimpan dalam bentuk *list of words*.

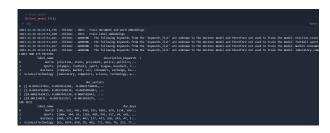


Fig. 2. Label beserta embedding vector dan dokumen terkait dari hasil  $\textit{training} \mod \text{Lbl2Vec}.$ 

Pemanggilan *function* .fit() pada gambar 2 digunakan untuk melakukan *training* model yang akan digunakan dalam

eksperimen ini. Pada proses *training*, terjadi beberapa hal penting sebagai berikut.

- Melakukan embedding vektor untuk label berdasarkan keyword deskripsi masing-masing label.
- 2) Melakukan embedding vektor untuk data training.
- 3) Menghitung *similarity* dari hasil vektor embedding setiap dokumen *training* dengan hasil vektor embedding label.
- 4) Mengategorikan setiap dokumen ke dalam label dengan *similarity* tertinggi.

Gambar 2 menampilkan dua hal, pertama nama label beserta keyword deskripsi dan vektor embedding label dan kedua nama label beserta index data yang termasuk ke dikategorikan sebagai label tersebut.

```
model_docs_lbl_similarities = lbl2vec_model_predict_model_docs()
2023-11-24 14:23:13,796 - Lbl2Vec - INFO - Get document embeddings from model
2023-11-24 14:23:13,798 - Lbl2Vec - INFO - Calculate document<->label similarities
             most similar label highest similarity score
                                                  0.000814
                          World
                                                  0.043607 0.037535
                         Sports
                                                            0.064150
        1196
1198
                                                  0.146370
                                                            0.052269
        1198
               -0.067967
               0.167956
      0.011518
                                    0.015231
```

Fig. 3. Hasil perhitungan similarity antara dokumen training dengan label.

Gambar 3 menampilkan hasil pengkategorian label untuk setiap dokumen *training* beserta skor *similarity* dengan semua label yang ada. Hasil ini didapatkan dari proses perhitungan *cosine similarity* dari vektor embedding dokumen dengan vektor embedding label.

```
# product idelically serves of one loci decounts (they were not used during (these: training)
now_dex_bl_stall artise = hithree_model_predict_new_dexit_aged_deco-ag_foll_corpos["faged_deco"][ag_foll_corpos["fafa_set_hype"]="text"])

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = 100 - Calculate document esteddings

## 203-11-24 1400103,pps = lbDrev = lb
```

Fig. 4. Hasil perhitungan similarity antara dokumen testing dengan label.

Gambar 4 menampilkan hasil pengkategorian label untuk setiap dokumen *testing* beserta skor *similarity* dengan semua label yang ada. Hasil ini didapatkan melalui proses yang sama dengan proses untuk mendapatkan hasil pada gambar 3, yaitu dengan menggunakan *cosine similarity*.

# III. HASIL UJI COBA DAN ANALISIS

#### A. Metrik Evaluasi

Metrik evaluasi yang digunakan pada eksperimen ini adalah F1-score. F1-score dihitung dengan cara mempertimbangkan precision dan recall secara imbang[2].

$$F1score = \frac{2 \times precision \times recall}{precision + recall} \tag{1}$$

## B. Dataset

Dataset yang digunakan dalam eksperimen ini adalah dataset AG's Corpus¹ dengan detail seperti pada pada tabel I. Dataset tersebut menggunakan bahasa Inggris. Pada eksperimen ini, dataset disimpan dalam bentuk csv. File train.csv yang berisi corpus untuk training dan test.csv yang berisi corpus untuk testing.

TABLE I
DETAIL DATASET[1]

Dataset	#Training documents	#Test documents	#Classes
AG's Corpus	120000	7600	4

# C. Experiment Settings

Eksperimen ini tidak membutuhkan spesifikasi *hardware* atau *software* khusus dalam menjalankan program. Saya menggunakan beberapa model dengan *hyperparameters* yang berbeda dalam melakukan eksperimen ini yang dijelaskan pada tabel II.

TABLE II
Experiments settings

Model	Hyperparameters	
	$similarity\_threshold = 0.30$	
A	$min\_num\_docs = 100$	
	epochs = 10	
	$similarity\_threshold = 0.50$	
В	$min\_num\_docs = 100$	
	epochs = 10	
	$similarity\_threshold = 0.30$	
C	$min\_num\_docs = 200$	
	epochs = 10	
	$similarity\_threshold = 0.30$	
D	$min\_num\_docs = 100$	
	epochs = 15	

<sup>&</sup>lt;sup>1</sup>https://github.com/mhjabreel/CharCnn\_Keras/tree/master/data/ag\_news\_csy

## D. Hasil

Hasil skor eksperimen untuk masing-masing model dapat dilihat pada tabel III. Model A (baseline model) memberikan hasil F1-score yang cukup baik, yaitu 0.8219 saat training dan 0.8151 saat testing. Kemudian, peningkatan nilai similarity\_threshold pada model B berdampak pada turunnya F1-score sebesar 0.0104 pada training dan 0.0116 pada testing. Sedangkan, peningkatan nilai min\_num\_docs pada model C memberikan hasil F1-score yang paling tinggi, baik pada saat training maupun testing dengan skor sebesar 0.8233 dan 0.8171. Terakhir, peningkatan jumlah epochs pada model D berdampak pada penurunan F1-score training sebesar 0.002, namun peningkatan F1-score testing sebesar 0.0012.

TABLE III HASIL EKSPERIMEN

	F1-score		
Model	Training	Testing	
A	0.8219	0.8151	
В	0.8115	0.8035	
С	0.8233	0.8171	
D	0.8199	0.8163	

## E. Diskusi

Pada skenario model B, peningkatan similarity\_threshold menyebabkan penurunan F1-score untuk kedua fase *training* dan *testing*. *Hyperparameter* ini digunakan sebagai *threshold/*batas minimum untuk mengategorikan suatu dokumen menjadi suatu label. Semakin tinggi nilai *hyperparameter* ini, maka semakin sulit untuk mengategorikan suatu dokumen. Oleh karena itu terjadi penurunan F1-score pada tahap *training* dan *testing*.

Pada skenario model C, peningkatan min\_num\_docs memberikan hasil F1-score tertinggi untuk kedua fase *training* dan *testing*. *Hyperparameter* ini merupakan jumlah minimal dokumen yang akan dipilih untuk menghitung *label embeddings*. Semakin tinggi nilai *hyperparameter* ini, maka model akan menghasilkan dokumen dengan informasi yang lebih signifikan dan kontekstual. Oleh karena itu, model akan lebih mudah dalam mengategorikan suatu dokumen dan menghasilkan F1-score yang lebih tinggi pada *training* dan *testing*.

Pada skenario model D, peningkatan jumlah *epochs* menyebabkan penurunan F1-score pada fase *training*, namun peningkatan pada fase *testing*. *Epochs* adalah *hyperparameter* yang digunakan untuk melatih seberapa lama model akan dilatih. Semakin tinggi nilai *epochs*, maka model akan semakin lama dilatih. Oleh karena itu, dengan meningkatnya jumlah *epochs* model akan lebih memahami pola dataset yang diberikan dan akan memberikan hasil evaluasi *testing* yang lebih bagus. Namun, perlu diwaspadai bahwa semakin tinggi nilai *epochs*, model akan semakin rawan menjadi *overfitting*.

Kelemahan model Lb12Vec adalah model tidak dapat belajar untuk menemukan *hyperparameter* sendiri yang bagus seperti pada *neural networks*. Melainkan, peneliti harus

menentukan *hyperparameter* sendiri yang didefinisikan di awal. Agar mendapatkan model dengan performa yang bagus, peneliti harus mencoba-coba berbagai kombinasi *hyperparameter* secara manual. Selain itu, kelemahan lain dari model ini adalah jika pada label terdapat kata/frasa yang belum pernah dipelajari/didapat dari dokumen *training*, maka kata/frasa tersebut akan di-*exclude* dan menyebabkan turunnya hasil F1-score.

## IV. KESIMPULAN DAN SARAN

# A. Kesimpulan

Pada eksperimen ini, penulis melakukan klasifikasi dokumen melalui pendekatan unsupervised dengan memanfaatkan model Lbl2Vec. Model ini bekerja berdasarkan vektor embedding dari dokumen sebagai unlabeled data dan label beserta deskripsi yang telah didefinisikan di awal. Pada eksperimen ini model C memberikan hasil F1-score terbaik, 0.8233 pada training serta 0.8171 pada *testing*. Model C menggunakan hyperparameter similarity threshold=0.30, min\_num\_docs=200, dan epochs=10. Hal menunjukkan bahwa model Lbl2Vec memberikan hasil F1-score yang cukup baik dalam melakukan klasifikasi teks secara unsupervised.

#### B. Saran

Eksperimen yang telah dilakukan oleh penulis masih sangat terbatas. Eksperimen ini hanya menggunakan 120.000 data training dan 7600 data testing yang berbahasa inggris. Oleh karena itu, perlu dilakukan eksperimen lebih lanjut pada dokumen berbahasa selain inggris. Hal ini dilakukan untuk mengetahui performa model Lbl2Vec pada dokumen yang berbahasa selain bahasa inggris. Selain itu, penulis menyarankan untuk melakukan eksperimen pada dokumen multibahasa.

#### ACKNOWLEDGMENT

Eksperimen ini didukung oleh Fakultas Ilmu Komputer, Universitas Indonesia. Penulis mengucapkan terima kasih banyak kepada Bu Ika Alfina selaku pengampu mata kuliah Pengolahan Bahasa Manusia pada Semester Gasal 2023/2024. Serta kepada Luthfi Balaka selaku asisten dosen yang membantu dalam proses penyelesaian eksperimen ini.

# REFERENCES

- [1] T. Schopf, D. Braun, and F. Matthes, "Lbl2vec: An embedding-based approach for unsupervised document retrieval on predefined topics," in *Proceedings of the 17th International Conference on Web Information Systems and Technologies*. SCITEPRESS Science and Technology Publications, 2021. [Online]. Available: http://dx.doi.org/10.5220/0010710300003058
- [2] M. Sokolova, N. Japkowicz, and S. Szpakowicz, "Beyond accuracy, f-score and roc: A family of discriminant measures for performance evaluation," vol. Vol. 4304, 01 2006, pp. 1015–1021.