

KLASIFIKASI BAKU LAPANGAN USAHA INDONESIA(KBLI) MENGUNAKAN *NATURAL LANGUAGE PROCESSING*

Lidya Elisabet Theogracia Silitonga¹, I Putu Gede Hendra Suputra², dan I Gede Santi Astawa³

ABSTRAK

Klasifikasi Baku Lapangan Usaha Indonesia(KBLI) merupakan klasifikasi baku kegiatan ekonomi yang terdapat di Indonesia. KBLI disusun untuk menyediakan satu set kerangka klasifikasi kegiatan ekonomi yang komprehensif di Indonesia agar dapat digunakan untuk penyeragaman, pengumpulan, pengolahan, penyajian, dan analisis data statistik menurut kegiatan ekonomi serta untuk mempelajari keadaan atau perilaku ekonomi menurut kegiatan ekonomi. Badan Pusat Statistik Provinsi Bali memiliki data baku lapangan usaha di Indonesia, akan tetapi data tersebut masih berbentuk data mentah dan belum diklasifikasikan secara menyeluruh. Implementasikan *Natural Language Processing* dalam proses klasifikasi data baku lapangan usaha Indonesia yang bertujuan untuk meningkatkan akurasi model dan meningkatkan performa klasifikasi yang akurat.

Kata kunci : klasifikasi baku lapangan usaha Indonesia

ABSTRACT

The Standard Classification of Business Fields in Indonesia (KBLI) is a standard classification of economic activities in Indonesia. The KBLI was developed to provide a comprehensive set of economic activity classification frameworks in Indonesia so that they can be used for standardizing, collecting, processing, presenting, and analyzing statistical data according to economic activity and for studying economic conditions or behavior according to economic activity. The Central Bureau of Statistics for the Province of Bali has standard data on business fields in Indonesia, but this data is still in the form of raw data and has not been thoroughly classified. Implement transfer learning in the process of classifying raw data for business fields in Indonesia which aims to improve model accuracy and improve accurate classification performance.

Keywords: Indonesian Standard Industrial Classification

1. PENDAHULUAN

Klasifikasi Baku Lapangan Usaha di Indonesia(KBLI) merupakan klasifikasi baku kegiatan ekonomi yang terdapat di Indonesia. KBLI disusun untuk menyediakan satu set kerangka klasifikasi kegiatan ekonomi yang komprehensif di Indonesia agar dapat digunakan untuk penyeragaman, pengumpulan, pengolahan, penyajian, dan analisis data statistik menurut kegiatan ekonomi serta untuk mempelajari keadaan atau perilaku ekonomi menurut kegiatan ekonomi. Badan Pusat Statistik Provinsi Bali memiliki data baku lapangan usaha di Indonesia, akan tetapi data tersebut masih berbentuk data mentah dan belum diklasifikasikan secara menyeluruh. Oleh karena itu, penulis mengimplementasikan *Natural Language Processing(NLP)* dalam

¹ Informatika, Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Udayana

² Informatika, Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Udayana

³ Informatika, Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Udayana

Submitted: 8 Oktober 2023

Revised: 2 November 2023

Accepted: 3 November 2023

Klasifikasi Baku Lapangan Usaha Indonesia(KBLI) Menggunakan Natural Language Processing(NLP)

proses klasifikasi data baku lapangan usaha di Indonesia yang bertujuan untuk meningkatkan akurasi model dan meningkatkan performa klasifikasi yang akurat.

2. METODE PELAKSANAAN

Adapun tahapan proses klasifikasi baku lapangan usaha di Indonesia menggunakan *Natural Language Processing*, yaitu

- **Pemahaman Dataset (*Understanding Dataset*)**
Pemahaman dataset yang dimiliki yang bertujuan untuk memperoleh informasi dari data tersebut.
- **Loading the Data**
Sebelum melakukan *Loading the Data*, kita harus menginstall terlebih dahulu modul dan mengimport *libraries* yang dibutuhkan. Setelah itu, melakukan *loaded dataset* yang mana data tersebut akan disimpan setekah dibaca dari file CSV menggunakan DataFrame(df).
- **Text processing**
Pada tahap ini, data tersebut akan diubah menjadi vektor dengan menggunakan TFIDF(Term Frequency-Inverse Document Frequency).
- **Mengeksplor Model Klasifikasi Teks**
Adapun model klasifikasi yang digunakan, yaitu Random Forest, Linear Support Vector Machine, Multinomial Naïve Bayes, dan Logistic Regression)
- **Membandingkan Performance Model Klasifikasi**
Pada tahap ini akan membandingkan ‘Mean Accuracy’ dan ‘Standar Deviation’ dari ke empat algoritma tersebut (model yang cocok untuk proses klasifikasi data tersebut).
- **Evaluasi Model Klasifikasi Teks**
Pada tahap ini akan dilakukan *train* terhadap model data yang dipilih sehingga kita dapat mengevaluasi dan memeriksa *performance* data yang tidak terlihat.
- **Prediksi**
Melakukan prediksi pada data yang tidak terlihat dan mengecek *performance* dari model.

3. HASIL DAN PEMBAHASAN

Berikut ini proses klasifikasi data baku lapangan usaha Indonesia menggunakan *Natural Language Processing*:

- **Loading the Data**
Pada tahap *loading the Data*, download terlebih dahulu data yang kita miliki dengan melakukan import dataset dari google drive, setelah itu install modul yang dibutuhkan dan kita juga dapat mengecek jumlah baris dan kolom dari dataset, seperti pada gambar 3.1. Pada gambar 3.2 dibuat data frame baru yang mana pada awalnya dataset tersebut memiliki 65535 baris dan kolom, lalu pada saat membuat data frame baru terdapat dua kolom yang dipilih, yaitu B1R15A_KOD(Kode Usaha) dan B1R15A(Jenis Usaha) untuk data frame baru serta menghapus *missing value*(NaN) yang terdapat pada baris dataset.

```
[ ] from google.colab import drive
drive.mount('/content/drive')
Mounted at /content/drive

[ ] import os
import pandas as pd
import numpy as np
from scipy.stats import randint
import seaborn as sns # used for plot interactive graph.
import matplotlib.pyplot as plt
import seaborn as sns
from io import StringIO
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.feature_selection import chi2
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.naive_bayes import MultinomialNB
from sklearn.linear_model import LogisticRegression
from sklearn.svm import LinearSVC
from sklearn.model_selection import cross_val_score
from sklearn.metrics import confusion_matrix
from sklearn import metrics

[ ] df = pd.read_csv('/content/drive/MyDrive/NBLI_BESC.csv')
print(df.shape)

(65535, 5)
```

Gambar 3.1. Install modul dan libraries & loaded dataset

```
[ ] # Create a new dataframe with two columns
df1 = df[['B1R15A_KOD', 'B1R15A']].copy()
# Remove missing values (NaN)
df1 = df1[pd.notnull(df1['B1R15A'])]
# Renaming second column for a simpler name
df1.columns = ['B1R15A_KOD', 'B1R15A']
print(df1.shape)
df1.head(3).T

(65531, 2)

B1R15A_KOD      0      1      2
B1R15A      C      G      C
B1R15A  BUYING AND SALE OF CEREMONY FACILITIES  SELLING SNACK FOOD  CARVING A STANDARD FROM WOOD
```

Gambar 3.2. Membuat data frame baru

Setelah membuat data frame baru dan menghapus *missing value*(NaN), pada gambar 3.3 dilakukan pengecekan kode apa saja yang terdapat di B1R15A_KOD(Kode Usaha). Lalu pada gambar 3.4, untuk mempermudah proses *training*, maka akan dilakukan beberapa perubahan pada nama kategori, yaitu B1R15A(Kode Usaha).

```
[ ] pd.DataFrame(df1.B1R15A_KOD.unique()).values

array([[ 'C',
       ['G'],
       ['I'],
       ['S'],
       ['M'],
       ['P'],
       ['K'],
       ['G'],
       ['P'],
       ['W'],
       ['P'],
       ['W'],
       ['Q'],
       ['E'],
       ['E'],
       ['E'],
       ['E'],
       ['E'],
       ['E']], dtype=object)

[ ] df2 = df1.sample(10000, random_state=1).copy()
```

Gambar 3.3. Melakukan pengecekan nilai unik

```
# Renaming categories
df2.rename(columns={'B1R15A_KOD':
                  ['C', 'I', 'S', 'M', 'P', 'K', 'G', 'W', 'Q', 'E'],
                  inplace=True)
pd.DataFrame(df2.B1R15A_KOD.unique())
```

```
0  G
1  I
2  P
3  G
4  S
5  Q
6  J
7  K
8  L
9  E
10 H
11 N
12 R
13 F
```

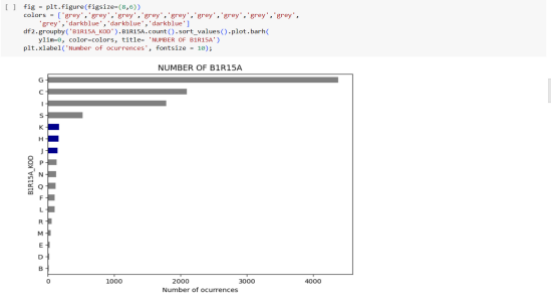
Gambar 3.4. Melakukan perubahan nama Kategori

Pada gambar 3.5 dilakukan pemetaan pada masing-masing kategori ke sebuah angka sehingga model dapat memahami dengan baik dan akan disimpan ke dalam kolom baru dengan nama 'category_id' yang mana masing-masing dari kategori direpresentasikan secara numerik. Pada gambar 3.6 dilakukan visualisasi dari data tersebut untuk melihat jumlah Jenis Usaha(B1R15A) yang terdapat pada tiap kategori(B1R15A_KOD).

```
[ ] # Create a new column 'category_id' with encoded categories
df2['category_id'] = df2['B1R15A_KOD'].factorize()[0]
category_id_df = df2[['B1R15A_KOD', 'category_id']].drop_duplicates()
# Dictionaries for future use
category_to_id = dict(category_id_df.values)
id_to_category = dict(category_id_df[['category_id', 'B1R15A_KOD']].values)
# New dataframe
df2.head()
```

B1R15A_KOD	B1R15A	category_id
12294	G SELLING SHOES	0
56925	G SELLING FOODS AND TOBACCO	0
43007	G SELLING COURSE BOOKS	0
35330	G SELLING SIDE FOOD	0
1131	I SELLING FOOD IN SHOPS	1

Gambar 3.5. Map masing-masing kategori



Gambar 3.6. Visualisasi data

- Text processing**

Pada tahap *text processing* dan pada gambar 3.7 teks diubah menjadi vektor menggunakan Term Frequency-Inverse Document Frequency (TFIDF) yang mana *TFIDFVectorizer* digunakan dengan tiga parameter, yaitu *min_df* digunakan untuk menghapus kata-kata yang muncul dalam jumlah file kurang dari 'min_df', *sublinear_tf* digunakan untuk menskalakan frekuensi dalam skala logaritmik, dan *stop_words* untuk menghapus kata-kata yang tidak memiliki makna.

```
[ ] tfidf = TfidfVectorizer(sublinear_tf=True, min_df=0,
                          ngram_range=(1, 2),
                          stop_words='english')
# fit transform each document into a vector
features = tfidf.fit_transform(df2.B1R15A).toarray()
labels = df2.category_id

print("Each of the 10000 B1R15A is represented by 582 Features (TF-IDF score of unigrams and bigrams)" % features.shape)

# Finding the three most correlated terms with each of the product categories
for Product, category_id in sorted(category_id.items()):
    features_chi2 = chi2(features, labels == category_id)
    indices = np.argsort(features_chi2)[-3:]
    feature_names = np.array([features_chi2[feature_names_out][i] for i in indices])
    unigrams = [w for w in feature_names if len(w.split()) == 1]
    bigrams = [w for w in feature_names if len(w.split()) == 2]
    print("\n %s - %s" % (Product, unigrams))
    print(" * most correlated unigrams are: %s" % ", ".join(unigrams[-3:]))
    print(" * most correlated bigrams are: %s" % ", ".join(bigrams[-3:]))

# Product: 'G'
# most correlated unigrams are: sell, need, whom
# most correlated bigrams are: selling shoes, buying sell, new building
# Product: 'I'
# most correlated unigrams are: processing, industry, shop
# most correlated bigrams are: sewing clothes, garment facilities, processing industry
# Product: 'S'
# most correlated unigrams are: electric, electricity, cubes
# most correlated bigrams are: electric poles, new ice cubes
# Product: 'M'
# most correlated unigrams are: water, water, collector
# most correlated bigrams are: drinking water, goods collector, used goods
# Product: 'P'
# most correlated unigrams are: container, building, construction
# most correlated bigrams are: building building, building building, building construction
```

Gambar 3.7. Mengubah text menjadi vektor menggunakan TFIDF dan Menemukan korelasi

- **Exploring Multi-Classification Models**

Pada gambar 3.8 data akan dibagi menjadi train dan test set yang mana data yang digunakan untuk training adalah 75% dan data yang digunakan untuk testing adalah 25%. Pada gambar 3.9, untuk membuat model klasifikasi terdapat beberapa algoritma yang digunakan, yaitu RandomForestClassifier, LinearSVC, MultinomialNB, dan LogisticRegression.

```
[ ] X = df2['BIRISA']
y = df2['BIRISA_KOD']
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                test_size=0.25,
                                                random_state = 0)
```

Gambar 3.8. Split data ke train dan test sets

```
[ ] models = [
    RandomForestClassifier(n_estimators=100, max_depth=5, random_state=0),
    LinearSVC(),
    MultinomialNB(),
    LogisticRegression(random_state=0),
]
# 5 Cross-validation
CV = 5
cv_df = pd.DataFrame(index=range(CV * len(models)))
entries = []
for model in models:
    model_name = model.__class__.__name__
    accuracies = cross_val_score(model, features, labels, scoring='accuracy', cv=CV)
    for fold_idx, accuracy in enumerate(accuracies):
        entries.append((model_name, fold_idx, accuracy))
cv_df = pd.DataFrame(entries, columns=['model_name', 'fold_idx', 'accuracy'])
```

Gambar 3.9. Menghitung dan Membandingkan performance model

- **Compare Text Classification Model Performance**

Setelah melakukan percobaan terhadap empat algoritma tersebut, akurasi dari algoritma LinearSVC lebih tinggi dari algoritma lainnya yang mana algoritma tersebut digunakan untuk model klasifikasi.



Gambar 3.10 Proses klasifikasi menggunakan model LinearSVC akurat.

- **Evaluation of Text Classification Model**

Pada gambar 3.11 model akan dilatih menggunakan 'LinearSVC' sehingga kita dapat mengevaluasi dan mengecek performance pada data yang tidak terlihat. Lalu pada gambar 3.12 akan dilakukan plot matriks untuk memeriksa *miss classified prediction*. Dan gambar 3.13 merupakan hasil dari plot *confusion matrix*.

```
[ ] X_train, X_test, y_train, y_test, indices_train, indices_test = train_test_split(features,
                                                                                labels,
                                                                                df2.index, test_size=0.25,
                                                                                random_state=1)

model = LinearSVC()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)

print('*****CLASSIFICATION METRICS*****')
print(metrics.classification_report(y_test, y_pred,
                                    target_names=df2['BIRISA_KOD'].unique()))

*****CLASSIFICATION METRICS*****
precision    recall  f1-score   support

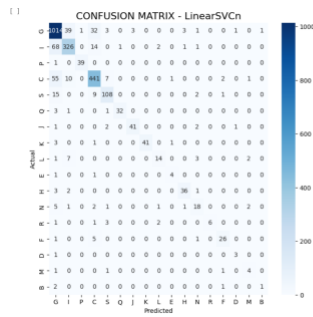
G   0.86   0.92   0.89   1008
I   0.84   0.79   0.82   413
P   0.97   0.97   0.97    68
C   0.87   0.85   0.86   517
S   0.86   0.88   0.87   1255
Q   0.97   0.86   0.91    17
J   0.93   0.87   0.90    47
K   1.00   0.89   0.94    86
L   0.74   0.52   0.61    27
E   0.67   0.67   0.67     6
H   0.86   0.86   0.87    42
N   0.62   0.58   0.60    31
R   1.00   0.46   0.63    13
F   0.84   0.79   0.81    13
D   0.80   0.75   0.77     4
M   0.44   0.57   0.50     7
B   0.50   0.25   0.33     4

accuracy   0.80   0.73   0.86   2500
macro avg  0.80   0.73   0.86   2500
weighted avg  0.86   0.86   0.86   2500
```

Gambar 3.11. Train Model Menggunakan Linear Support Vector Machine

```
conf_mat = confusion_matrix(y_test, y_pred)
fig, ax = plt.subplots(figsize=(8,8))
sns.heatmap(conf_mat, annot=True, cmap='Blues', fmt='d',
            xticklabels=category_id_df.BIRISA_KOD.values,
            yticklabels=category_id_df.BIRISA_KOD.values)
plt.xlabel('Actual')
plt.ylabel('Predicted')
plt.title('CONFUSION MATRIX - LinearSVC', size=10);
```

Gambar 3.12. Plot confusion matrix



Gambar 3.13. Confusion Matrix

- *Prediction*
 Pada gambar 3.14 dibuat prediksi pada data yang tidak terlihat dan mengecek performance dari model yang telah dibuat.

```
[ ] X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.25,
                                                    random_state = 0)

tfidf = TfidfVectorizer(sublinear_tf=True, min_df=5,
                        ngram_range=(1, 2),
                        stop_words='english')

fitted_vectorizer = tfidf.fit(X_train)
tfidf_vectorizer_vectors = fitted_vectorizer.transform(X_train)
model = LinearSVC().fit(tfidf_vectorizer_vectors, y_train)
```

Gambar 3.14. Membuat prediksi dan Mengecek performance model

Pada gambar 3.15, misal kita memasukkan jenis usaha(B1R15A) adalah “SELLING COFFEE” maka setelah *code* di *run*, maka akan ditampilkan kode usaha(B1R15A_KOD) adalah ‘I’.

```
[ ] B1R15A = "SELLING COFFEE"
    print(model.predict(fitted_vectorizer.transform([B1R15A])))

['I']
```

Gambar 3.15. Hasil Proses Klasifikasi

4. KESIMPULAN

Badan Pusat Statistik Provinsi Bali memiliki data baku lapangan usaha di Indonesia, akan tetapi data tersebut masih berbentuk data mentah dan belum diklasifikasikan secara menyeluruh. Implementasikan transfer learning dalam proses klasifikasi data baku lapangan usaha di Indonesia yang bertujuan untuk meningkatkan akurasi model, meningkatkan performa klasifikasi yang akurat.

DAFTAR PUSTAKA

Muthia, D. A. (2018). Komparasi Algoritma Klasifikasi Text Mining Untuk Analisis Sentimen Pada Review Restoran. *Jurnal PILAR Nusa Mandiri*, 14(1), 69-74.

Halaman ini sengaja dikosongkan