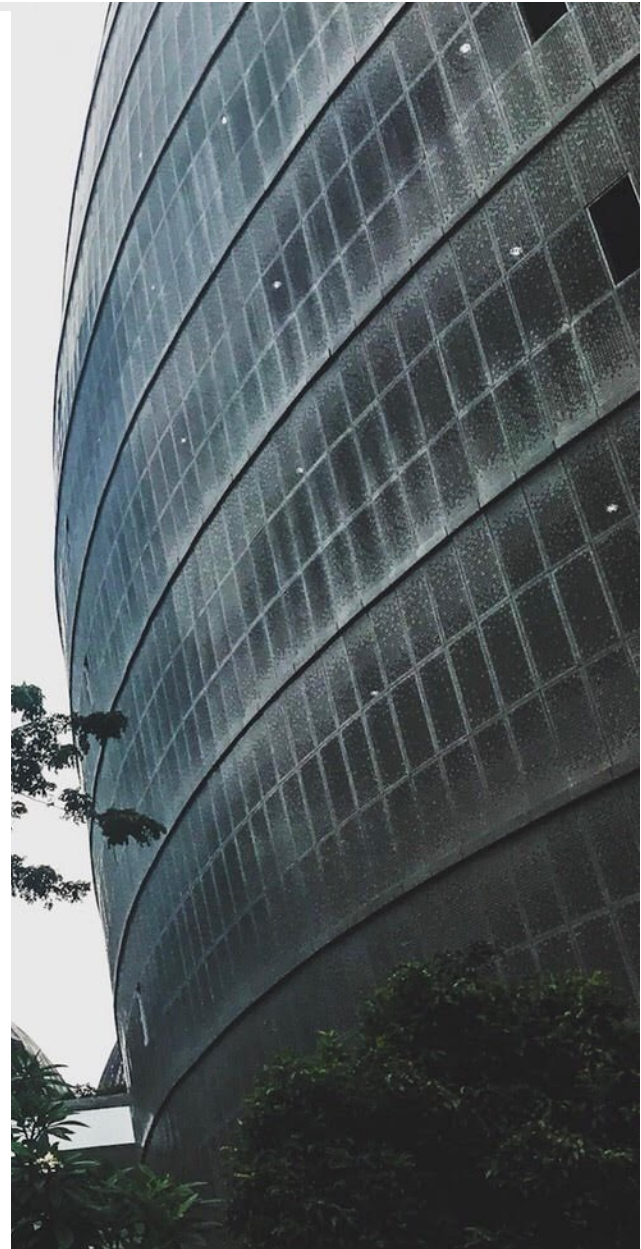


# MODUL PRAKTIKUM

IS411 – DATA MODELING  
PROGRAM SARJANA S1 SISTEM INFORMASI  
FAKULTAS TEKNIK DAN INFORMATIKA

---



---

**PROGRAM STUDI SISTEM INFORMASI  
FAKULTAS TEKNIK DAN INFORMATIKA  
UNIVERSITAS MULTIMEDIA NUSANTARA**

Gedung B Lantai 5, Kampus UMN  
Jl. Scientia Boulevard, Gading Serpong, Tangerang, Banten-15811 Indonesia  
Telp: +62-21.5422.0808 (ext. 1803), email: [ict.lab@umn.ac.id](mailto:ict.lab@umn.ac.id), web: [umn.ac.id](http://umn.ac.id)

# MODUL 8

## MACHINE LEARNING



### DESKRIPSI TEMA

1. Training and Testing Set on Machine Learning
2. Further Knowledge to Classification Algorithms.

### CAPAIAN PEMBELAJARAN MINGGUAN (SUB-CAPAIAN PEMBELAJARAN)

Students are able to implement machine learning techniques or methods (C3)

### PENUNJANG PRAKTIKUM

1. Anaconda Navigator
  2. Jupyter Notebook
- (+ Perlengkapan Apd/Alat Pelindung Diri Yang Harus Digunakan, Jika Ada)

### LANGKAH-LANGKAH PRAKTIKUM

#### Import Library

1. Lakukan langkah yang sama seperti minggu sebelumnya untuk meng-import library Numpy dan Pandas.
2. Import dataset menggunakan file “winequality-red.csv”.

```
1 dataset = pd.read_csv('winequality-red.csv', delimiter=';')
2 dataset
```

3. Lakukan pembagian data training dan testing.

```
1 from sklearn.model_selection import train_test_split
2
3 #membagi data ke X dan Y
4 redwine = dataset.copy()
5 Y = redwine['quality']
6 X = redwine.drop(columns = 'quality')
7
8 #membagi data training 80% dan testing 20%
9 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2, random_state=0)
```

#### Cross Validation

Pada praktikum sebelumnya telah dilakukan klasifikasi menggunakan algoritma LogisticRegression terhadap data wine, dimana setelah melakukan testing dan training tidak ada validasi terhadap hasil klasifikasi.

Sebagai perbandingan terhadap hasil dari klasifikasi, berikut adalah contoh klasifikasi LogisticRegression biasa.

Berikut ini adalah contoh data redwine dengan label yang sudah diberikan.

4. Pemodelan dengan algoritma LogisticRegression.

```
1 from sklearn.linear_model import LogisticRegression
2
3 logreg = LogisticRegression()
4 logreg.fit(X_train, y_train)
5 print('Accuracy of Logistic regression classifier on training set: {:.2f}'.format(logreg.score(X_train, y_train)))
6 print('Accuracy of Logistic regression classifier on test set: {:.2f}'.format(logreg.score(X_test, y_test)))
```

5. Selanjutnya adalah melakukan Cross Validation menggunakan library **cross\_val\_score**, sistem akan melakukan pengulangan pada data train yang akan menghasilkan akurasi yang berbeda.

```
1 from sklearn.model_selection import cross_val_score
2
3 scores = cross_val_score(logreg, X_train, y_train, cv=10)
```

```
1 print('Cross-Validation Accuracy Scores', scores)
2 scores = pd.Series(scores)
3 scores.min(), scores.mean(), scores.max()
```

6. Catatlah berapa hasil Cross-Validation score yang diperoleh.

## Natural Language Processing (NLP)

Natural language processing umumnya digunakan pada kumpulan teks, seperti yang dihasilkan / didapatkan dari Tweet, *post* media sosial, percakapan, *review* film, berita, dan sebagainya. Kumpulan teks disebut dengan **corpus**, atau dalam bentuk jamak disebut **corpora**.

Ada beberapa cara dan *tools* untuk melakukan NLP, pada praktikum kali ini kita akan menggunakan TextBlob. **TextBlob** adalah *object-oriented NLP text-processing library* yang dibangun pada NLTK dan *pattern NLP libraries*.

7. Membuat TextBlob.

```
1 #install TextBlob
2
3 #Cara 1: menggunakan Anaconda Prompt
4 #conda install -c conda-forge textblob
5
6 #Cara 2: menggunakan Jupyter Notebook
7 #!pip install textblob
8
9 !pip install textblob
```

```
1 #download the NLTK corpora used by TextBlob
2
3 import nltk
4 nltk.download('averaged_perceptron_tagger')
```

```
1 import nltk
2 nltk.download('brown')
```

```
1 #Create a TextBlob
2
3 from textblob import TextBlob
4
5 text = 'Today is a beautiful day. Tomorrow looks like bad weather.'
6
7 blob = TextBlob(text)
8
9 blob
```

8. Melakukan **Tokenisasi** teks ke dalam kalimat (*sentences*) dan kata (*words*).

```
1 #sentence property to get a List of Sentence objects
2
3 blob.sentences
```

```
1 #words property returns a WordList object containing a List of Word objects
2
3 blob.words
```

9. Memecah teks menjadi **Part-of-Speech Tagging**

*Part-of-Speech Tagging* memecah teks / kalimat ke dalam *tuples*, biasanya per suku kata.

```
1 #tags property returns a list of tuples, each containing a word and a string representing its part-of-speech tag
2
3 blob.tags
```

10. Kita juga dapat mengekstrak kalimat atau teks ke dalam frasa.

```
1 #noun_phrases property returns a WordList object containing a list of Word objects - one for each noun phrase in the text
2
3 blob.noun_phrases
```

11. Melakukan analisis sentimen (*sentiment analysis*) menggunakan **Sentiment Analyzer** dari TextBlob

```
1 #getting the sentiment of a TextBlob
2
3 blob.sentiment
```

```
1 #getting the polarity and subjectivity from the Sentiment Object
2
3 %precision 3
```

```
1 blob.sentiment.polarity
```

```
1 blob.sentiment.subjectivity
```

```
1 #getting the sentiment of a sentence
2 #get the sentiment at the individual sentence level
3
4 for sentence in blob.sentences:
5     print(sentence.sentiment)
6
```

## 12. Melakukan analisis sentiment menggunakan algoritma Naïve Bayes, dengan **NaiveBayesAnalyzer**

```
1 from textblob.sentiments import NaiveBayesAnalyzer
2
3 blob = TextBlob(text, analyzer=NaiveBayesAnalyzer())
4
5 blob
```

Kita bisa mencoba menggunakan dataset lain, misalnya dataset ulasan film berikut.

```
1 import nltk
2 nltk.download('movie_reviews')
```

```
1 blob.sentiment
```

```
1 #get the sentiment of each sentences
2
3 for sentence in blob.sentences:
4     print(sentence.sentiment)
```

## 13. *Inflection: Pluralization dan Singularization*

**Inflection** adalah bentuk berbeda dari kata yang sama, biasanya dalam bahasa Inggris seringkali terdapat dua kata yang berbeda karena terdapat bentuk tunggal (*singular*) dan jamak (*plural*), meskipun memiliki arti yang sama. Misalnya “person” dan “people” sama – sama memiliki arti “orang”, atau “run” dan “ran” sama – sama memiliki arti “lari”.

```
1 #Words and WordLists each support converting words to their singular or plural forms.
2
3 from textblob import Word
4
5 index = Word('index')
6 index.pluralize()
```

```
1 cacti = Word('cacti')
2 cacti.singularize()
```

```

1 from textblob import TextBlob
2
3 animals = TextBlob('dog cat fish bird').words
4 animals.pluralize()

```

#### 14. Spell Checking and Correction

Kita dapat melakukan pengecekan atau koreksi terhadap ejaan yang keliru.

```

1 #spellcheck method returns a list of tuples containing possible correct spellings and a confidence value
2
3 from textblob import Word
4
5 word = Word('theyr')
6 word.spellcheck()

```

```

1 #correct method chooses word with the highest confidence value
2
3 word.correct()

```

```

1 from textblob import Word
2
3 sentence = TextBlob('Ths sentence has misspelled wrds.')
4 sentence.correct()

```

#### 15. Normalization: Stemming and Lemmatization

**Stemming** dilakukan untuk menghilangkan awalan (*prefix*) atau akhiran (*suffix*) dari suatu kata sehingga hanya tersisa sebagai kata dasar. Hal ini perlu dilakukan untuk mengurangi risiko adanya kata yang redundansi.

```

1 #Stemming removes a prefix or suffix from a word leaving only a stem, which may or may not a real word.
2 #Lemmatization is similar, but factors in the word's part of speech and meaning and results in a real word.
3
4 from textblob import Word
5
6 word = Word('strawberries')
7 word.lemmatize()

```

```

1 word.stem()

```

#### 16. Menghitung Frekuensi Kata

Untuk melakukan perhitungan frekuensi kata, kita menggunakan dataset **RomeoAndJuliet.txt**.

```

1 #Various techniques for detecting similarity between documents rely on word frequencies.
2
3 #use Path class form the Python Standard Library's pathlib module
4
5 from pathlib import Path
6 from textblob import TextBlob
7
8 #Load the e-book RomeoAndJuliet.txt (untuk memudahkan, letakan file di path yang sama dengan tempat menyimpan ipynb ini)
9 blob = TextBlob(Path('RomeoAndJuliet.txt').read_text())

```

```

1 blob.word_counts['juliet']

```

```
1 blob.word_counts['romeo']
```

**Pertanyaan: Berapa jumlah kata "juliet" dan "romeo" yang Anda dapatkan?**

17. Mendapatkan Arti Kata, Sinonim, dan Antonim menggunakan **WordNet**

**WordNet** adalah database yang dibuat oleh Princeton University. TextBlob library menggunakan antarmuka library NLTK yaitu WordNet, untuk mendapatkan arti kata, sinonim, dan antonim dari suatu kata.

```
1 #getting definitions
2
3 from textblob import Word
4
5 happy = Word('happy')
6
7 #definitions property returns a list of all the word's definitions in the WordNet database
8
9 happy.definitions
```

```
1 #getting synonyms
2
3 happy.synsets
```

18. Menghapus *Stop Word*

**Stop word** adalah kata umum pada teks yang biasanya dihilangkan sebelum melakukan analisis, karena dianggap tidak memberikan arti/makna penting untuk analisis.

```
1 #download NLTK's stop-words lists
2
3 import nltk
4
5 nltk.download('stopwords')
```

```
1 #load the english stop words list
2
3 from nltk.corpus import stopwords
4
5 stops = stopwords.words('english')
```

```
1 #create a TextBlob from wich we'll remove stop words
2
3 from textblob import TextBlob
4
5 blob = TextBlob('Today is a beautiful day.')
6
7 [word for word in blob.words if word not in stops]
```

19. N-Grams



**N-Grams** adalah urutan n item teks, seperti huruf dalam kata atau kata dalam kalimat. Dalam *Natural Language Processing*, n-gram dapat digunakan untuk mengidentifikasi huruf atau kata yang sering muncul berdekatan satu sama lain. Untuk input pengguna berbasis teks, ini dapat membantu memprediksi huruf atau kata berikutnya yang akan diketik pengguna.

```
1 from textblob import TextBlob
2
3 text = 'Today is a beautiful day. Tomorrow looks like bad weather.'
4
5 blob = TextBlob(text)
6 blob.ngrams()
```

```
1 #n-grams consisting of 5 words
2
3 blob.ngrams(n=5)
```

## 20. Visualisasi Frekuensi Kata dengan Bar Charts dan Word Clouds

Untuk membuat visualisasi jumlah/frekuensi suatu kata dari dataset, kita bisa menggunakan Pandas.

### a. Menggunakan Bar Charts

```
1 #load the data
2
3 from pathlib import Path
4 from textblob import TextBlob
5
6 blob = TextBlob(Path('RomeoAndJuliet.txt').read_text())
```

```
1 #load the NLTK stopwords
2
3 from nltk.corpus import stopwords
4
5 stop_words = stopwords.words('english')
```

```
1 #getting the word frequencies
2
3 items = blob.word_counts.items()
```

```
1 #eliminating the stop words
2
3 items = [item for item in items if item[0] not in stop_words]
4 #item[0] gets the word from each tuple so we can check the whether it's in stop_words
```

```
1 #sorting the words by frequency
2
3 from operator import itemgetter
4
5 sorted_items = sorted(items, key=itemgetter(1), reverse=True)
6
7 #As sorted orders item's elements,
8 #it accesses the element at index 1 in each tuple via the expression itemgetter(1)
9 #The reverse=True keyword argument indicates that the tuples should be sorted in descending order
```



```

1 #getting the top 20 words
2
3 top20 = sorted_items[1:21]

```

```

1 #convert top20 to a DataFrame
2
3 import pandas as pd
4
5 df = pd.DataFrame(top20, columns=['word', 'count'])
6 df

```

```

1 #visualizing the DataFrame
2
3 axes = df.plot.bar(x='word', y='count', legend=False)

```

#### b. Menggunakan Word Clouds

```

1 #installing the wordcloud module
2
3 #Cara 1: menggunakan Anaconda Prompt
4 #conda install -c conda-forge wordcloud
5
6 #Cara 2: via Jupyter Notebook
7 #!pip install wordcloud
8
9 !pip install wordcloud

```

```

1 #loading the text
2
3 from pathlib import Path
4
5 text = Path('RomeoAndJuliet.txt').read_text()

```

```

1 #configuring the WordCloud object
2
3 from wordcloud import WordCloud
4
5 wordcloud = WordCloud(colormap='prism', background_color='white')

```

```

1 #generating the Word Cloud
2
3 wordcloud = wordcloud.generate(text)

```

```

1 #saving the Word Cloud as an image file
2
3 wordcloud = wordcloud.to_file('RomeoAndJulietWordCloud.png')

```

#### Challenge:

1. Ambil kata-kata dari satu dokumen sampel (Shakespeare-hamlet.txt) di corpus Gutenberg dan tampilkan.
2. Tampilkan 500 kata pertama dari teks corpus tersebut dalam sebuah paragraf.
3. Lakukan proses Tokenization, Stemming, Lemmatization dan POS Tags pada paragraf tersebut.

## PENGUMPULAN

1. File yang dikumpulkan terdiri dari:
  - a. File project (.ipynb)
  - b. File PDF berisi screenshot output dan jawaban (jika ada pertanyaan)
2. File di-compress (.zip) dan diberi nama **KODEMATAKULIAH\_KELAS\_NIM\_NAMA\_WEEK KE-XX.zip**  
(contoh: IS5411\_A\_13110310017\_Monika Evelin Johan\_Week-01.zip).

## REFERENSI

Deitel, P., & Deitel, H. (2020). Intro to Python for Computer Science and Data Science. Pearson Education.