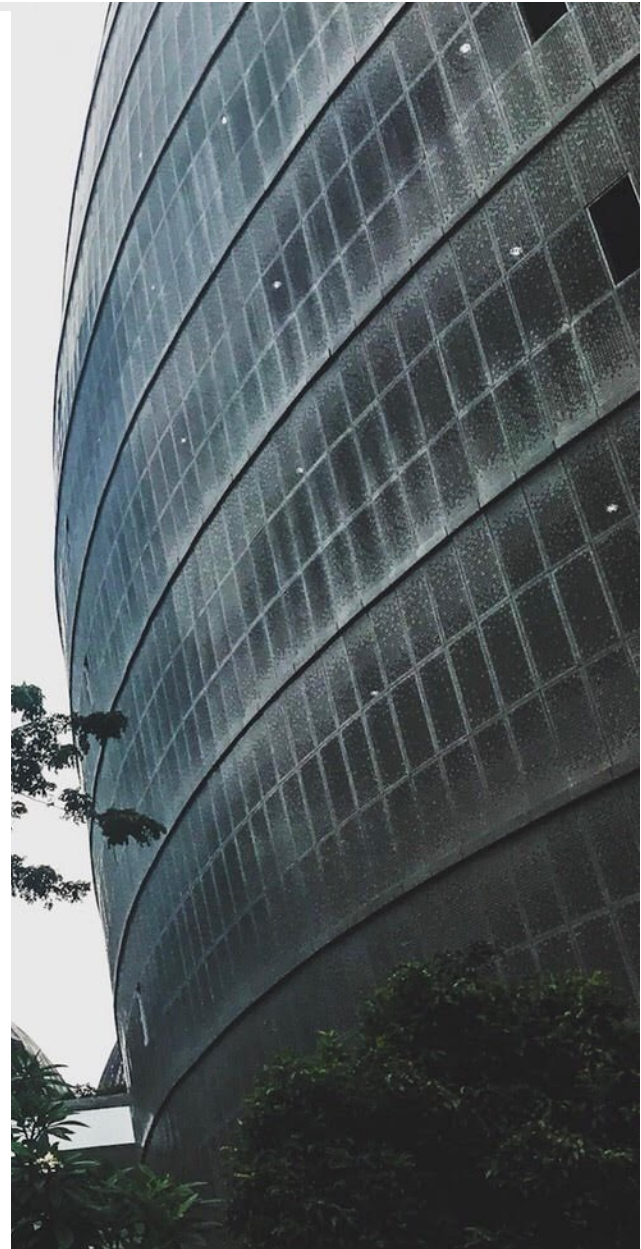


MODUL PRAKTIKUM

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, web: umn.ac.id

MODUL 9

CLASSIFICATION



DESKRIPSI TEMA

1. Decision Tree Algorithm
2. Random Forest Algorithm

CAPAIAN PEMBELAJARAN MINGGUAN (SUB-CAPAIAN PEMBELAJARAN)

Students are able to illustrate data modeling problems using classification techniques (C4)

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 untuk mengimport library Numpy, Pandas, Matplotlib, dan Seaborn
2. Import dataset menggunakan file "titanic_train.csv".

```
1 titanic = pd.read_csv('titanic_train.csv', delimiter=";")
2 titanic
```

Note: gunakan delimiter sesuai settingan pada komputer Anda (bisa ";" atau ",").

Data Exploration

3. Melihat deskripsi data

```
1 titanic.info()
```

Perhatikan bahwa terdapat beberapa atribut yang memiliki data null. Untuk mendapatkan hasil yang baik kita perlu mencoba memperbaiki data tersebut.

4. Mengisi Data Null

```

1 #Kolom Age
2 #mengisi data null pada kolom Age dengan rata - rata usia
3
4 rata_umur = titanic['Age'].mean()
5 titanic['Age'] = titanic['Age'].fillna(rata_umur)
6
7 #memeriksa apakah masih ada yang null pada kolom Age
8 titanic['Age'].isna().sum()

```

```

1 #Kolom Embarked
2 #mengisi data null pada kolom Embarked dengan nilai yang paling banyak muncul (modus)
3
4 modus = titanic['Embarked'].mode()[0]
5 titanic['Embarked'] = titanic['Embarked'].fillna(modus)
6 #memeriksa apakah masih ada yang null pada kolom Embarked
7 titanic['Embarked'].isna().sum()

```

5. Periksa apakah semua data sudah tidak ada yang null.

```

1 titanic.info()

```

Selain data null, terdapat beberapa atribut yang tidak akan digunakan sehingga perlu dihapus. Di samping itu terdapat kolom "Sex" dan kolom "Embarked" yang perlu diubah formatnya menjadi integer agar dapat digunakan pada modeling.

6. Data Cleansing

```

1 #Menghapus kolom Cabin
2
3 titanic.drop('Cabin', axis=1, inplace=True)
4 #axis = 1 menandakan kolom, axis = 0 menandakan baris

```

```

1 #membuang kolom Ticked
2 titanic.drop('Ticket', axis=1, inplace=True)

```

```

1 #menghapus kolom yang tidak digunakan
2 titanic.drop('Name', axis=1, inplace=True)

```

7. Konversi data

```

1 #mengkonver data kolom Sex, male = 0; female = 1
2
3 titanic['Sex'] = titanic['Sex'].replace(['male', 'female'], [0, 1])
4 titanic.head()

```

```

1 #mengkonver data kolom Embarked, S = 1, C = 2, Q = 3
2
3 titanic['Embarked'] = titanic['Embarked'].replace(['S', 'C', 'Q'], [1, 2, 3])
4 titanic.head()

```

```

1 titanic['Fare'] = titanic['Fare'].str.replace(".", "")

```

```

1 titanic.head()

```

Periksa apakah semua data sudah siap untuk digunakan.

Mencari Korelasi Data

8. Kita bisa melakukan analisis korelasi untuk mengetahui korelasi antar data sehingga bisa memilih hanya atribut yang paling berkorelasi dengan data target.

```

1 #untuk mencari tahu variabel / feature yang paling mempengaruhi target (Survived)
2
3 f, ax = plt.subplots(figsize=(10,10))
4
5 sns.heatmap(titanic.astype(float).corr(), linewidth=0.25, vmax=1.0, square=True, linecolor='black', annot=True)

```

Analisislah hasil dari Pearson Correlation tersebut. Mana saja atribut / feature yang menurutmu paling memiliki korelasi pada data target “Survived”?

Membangun Model

9. Menentukan nilai X dan Y

```

1 #membagi data X dan Y (x = variabel independen; y = variabel dependen/target)
2
3 x = titanic[['Sex', 'Parch', 'Fare', 'Embarked']]
4 #x hanya menggunakan data yang memiliki korelasi tinggi dengan Survived
5 y = titanic['Survived']

```

10. Membagi data training dan data testing

```

1 #membagi data training dan testing
2
3 from sklearn.model_selection import train_test_split
4
5 train_x, test_x, train_y, test_y = train_test_split(x, y, random_state = 0)
6 print(train_x.shape)
7 print(test_x.shape)
8 print(train_y.shape)
9 print(test_y.shape)

```

11. Modeling dengan algoritma Decision Tree

```

1 #Modeling dengan Decision Tree
2
3 from sklearn.tree import DecisionTreeClassifier
4
5 dt = DecisionTreeClassifier()
6 dt.fit(train_x, train_y)
7
8 #mengevaluasi hasil prediksi
9 y_prediksi3 = dt.predict(test_x)
10 akurasi_dt = dt.score(test_x, test_y)
11 akurasi_dt

```

12. Menampilkan struktur Decision Tree

```

1 #menampilkan struktur Decision Tree
2
3 from sklearn import tree
4
5 feature_names = x.columns[:4]
6 target_names = y.unique().tolist()
7
8 target_names = [str(x) for x in target_names]
9
10 fig = plt.figure(figsize=(25,20))
11 _ = tree.plot_tree(dt,
12                   feature_names = feature_names,
13                   class_names = target_names,
14                   filled = True)

```

13. Modeling dengan algoritma Random Forest

```

1 #modeling dengan Random Forest
2
3 from sklearn.ensemble import RandomForestClassifier
4
5 rf = RandomForestClassifier(n_estimators = 10, random_state = 42)
6 rf.fit(train_x, train_y)
7
8 #mengevaluasi hasil prediksi
9 y_prediksi4 = rf.predict(test_x)
10 akurasi_rf = rf.score(test_x, test_y)
11 akurasi_rf

```

Implementasi Machine Learning

Menggunakan dataset **titanic_test.csv** yang belum ada labelnya (tidak ada kolom Survived).

14. Lakukan langkah yang sama mulai dari import data, data cleansing, sampai menentukan data X.

Di sini kita tidak menentukan data Y karena tidak ada kolom target, yang akan dilakukan adalah prediksi.

```
1 titanic_test = pd.read_csv('titanic_test.csv', delimiter=';')
2 titanic_test.head()
```

```
1 titanic_test.info()
```

- Periksa data apa saja yang memiliki Null. Lakukan langkah yang sama seperti di atas untuk mengisi data tersebut. Untuk kolom Fare bisa mengisi data Null dengan nilai 0.
- Jangan lupa untuk mengubah format data Sex dan Embarked.

15. Menentukan data X.

```
1 #menentukan variabel X dan Y
2
3 x1 = titanic_test[['Sex', 'Parch', 'Fare', 'Embarked']]
```

16. Implementasi menggunakan Random Forest untuk prediksi

```
1 #implementasi Random Forest pada dataset untuk melakukan prediksi
2
3 titanic_prediction = rf.predict(x1)
4 titanic_prediction
```

17. Buat ke dataframe

```
1 #buat ke dataframe
2 df_titanic_prediction = pd.DataFrame({"Predicted Survived" : titanic_prediction})
3 df_titanic_prediction
```

18. Menggabungkan dataset dengan hasil prediksi

```
1 #menggabungkan hasil prediksi dengan dataset
2
3 result = pd.concat([titanic_test, df_titanic_prediction], ignore_index = False, axis = 1)
4 result
```

19. Menyimpan hasil yang baru pada file CSV

```
1 #menyimpan hasil yang baru ke dalam file csv
2 export_result2 = result.to_csv('Hasil Prediksi Titanic CSV.csv')
3 print(export_result2)
```

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.