Nama Dosen    : Teguh Iman Hermanto, M.Kom
Mata Kuliah   : Machine Learning 1
Pembahasan    : Naïve Bayes
Pokok Pemb    : - Konsep Algoritma Naïve Bayes
              - Membangun Model Naïve Bayes
              - Simulasi Algoritma Naïve Bayes
              - Evaluasi Algoritma Naïve Bayes
              - Aplikasi menggunakan algoritma Naïve Bayes
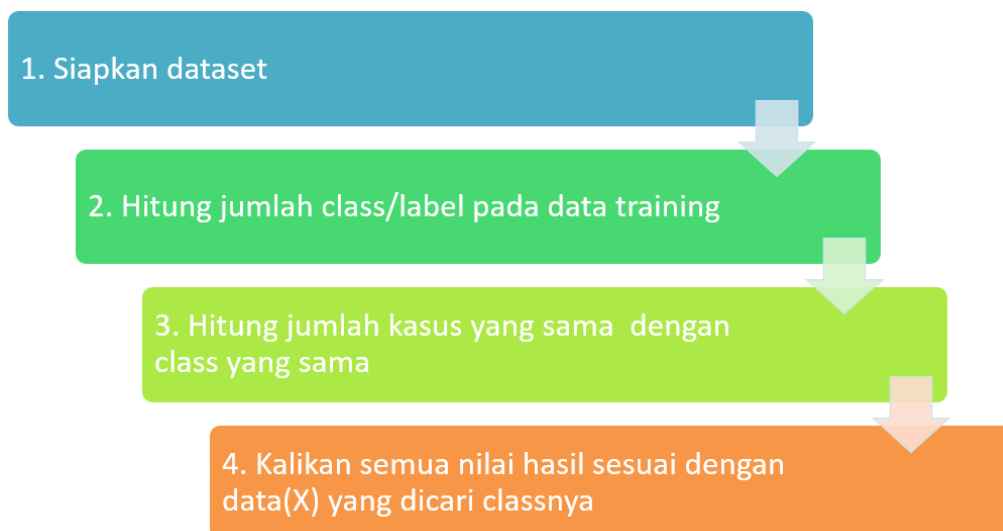
## 1. Konsep Algoritma Naïve Bayes

Naïve Bayes Classifier merupakan sebuah metoda klasifikasi yang berakar pada teorema Bayes . Metode pengklasifikasian dg menggunakan metode probabilitas dan statistik yg dikemukakan oleh ilmuwan Inggris Thomas Bayes , yaitu memprediksi peluang di masa depan berdasarkan pengalaman di masa sebelumnya sehingga dikenal sebagai Teorema Bayes . Ciri utama dr Naïve Bayes Classifier ini adalah asumsi yg sangat kuat (naïf) akan independensi dari masing-masing kondisi / kejadian.

Model Algoritma Naïve Bayes

$$P(H|X) = \frac{P(X|H)P(H)}{P(X)} = P(X|H) \times P(H)/P(X)$$

x        = Data dari class scoring yang belum diketahui
H        = Hipotesis data X yang merupakan suatu class yang lebih spesifik
P(H|X) = Probabilitas hipotesis H berdasarkan kondisi X
P(H)    = Probabilitas hipotesis H
P(X|H) = Probabilitas hipotesis X berdasarkan kondisi H
P(X)    = Probabilitas X

Tahapan Algoritma Naïve Bayes

1. Siapkan dataset

2. Hitung jumlah class/label pada data training

3. Hitung jumlah kasus yang sama  dengan class yang sama

4. Kalikan semua nilai hasil sesuai dengan data(X) yang dicari classnya

a. **Siapkan Dataset**
- Data training buys_computer dataset

| Age | Income | Student | Credit Rating | Buys_Computer |
|---|---|---|---|---|
| youth | high | no | fair | no |
| youth | high | no | excellent | no |
| middle_aged | high | no | fair | yes |
| senior | medium | no | fair | yes |
| senior | low | yes | fair | yes |
| senior | low | yes | excellent | no |
| middle_aged | low | yes | excellent | yes |
| youth | medium | no | fair | no |
| youth | low | yes | fair | yes |
| senior | medium | yes | fair | yes |
| youth | medium | yes | excellent | yes |
| middle_aged | medium | no | excellent | yes |
| middle_aged | high | yes | fair | yes |
| senior | medium | no | excellent | no |

- Data testing

| Age | Income | Student | Credit Rating | Buys_Computer |
|---|---|---|---|---|
| Senior | high | no | fair | ? |

Apakah calon pembeli akan membeli komputer?

b. **Hitung jumlah class/label pada data training**

Class 1 (C1)→ Buys_Computer yes → 9

Class 2 (C2) → Buys_Computer no → 5

Total Record: 14

Maka:

P(C1) = 9/14 = 0.642857143

P(C2) = 5/14 = 0.357142857

c. Hitung jumlah kasus yang sama dengan class yang sama

Contoh penghitungan $P(X|C_i)$ pada kasus Income dimana

i=1 dan i=2

Maka :

P(Income="high"|Buys_Computer="yes") = 2/9 = 0.22222222

P(Income="high"|Buys_Computer="no") = 2/5 = 0.4

**d. Kalikan semua nilai hasil sesuai dengan data(X) yang dicari classnya**

Jika semua data aribut pada data training dihitung, maka:

| Atribut | Parameter | yes | no |
|---|---|---|---|
| Age | youth | 0.222222222 | 0.6 |
| Age | middle_aged | 0.444444444 | 0 |
| **Age** | **senior** | **0.333333333** | **0.4** |
| Income | low | 0.333333333 | 0.2 |
| Income | medium | 0.444444444 | 0.4 |
| **Income** | **high** | **0.222222222** | **0.4** |
| Student | yes | 0.666666667 | 0.2 |
| **Student** | **no** | **0.333333333** | **0.8** |
| **Credit Rating** | **fair** | **0.666666667** | **0.4** |
| Credit Rating | excellent | 0.333333333 | 0.6 |

P(X|Buys_Computer="yes") = 0.016460905

P(X|Buys_Computer="no")  = 0.0512

P(X|Buys_Computer="yes")*P(C1)  = 0.010582011

P(X|Buys_Computer="no")*P(C2)   = 0.018285714

Nilai "no" lebih besar daripada nilai "yes". Sehingga pembeli dengan atribut X tidak membeli komputer

**2. Membangun Model Naïve Bayes**

```python
import pandas as pd
import numpy as np

from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import CategoricalNB
from sklearn.metrics import accuracy_score, confusion_matrix

import matplotlib.pyplot as plt
import seaborn as sns
```

```
1  df = pd.read_csv('Diabetes1.csv')
```

|   | Age | BMI | Blood_Sugar_Level | Family_History | Diet | Diabetes |
|---|-----|------|-------------------|----------------|----------|----------|
| 0 | 45 | 28.7 | 135 | Yes | Poor | Yes |
| 1 | 50 | 31.2 | 145 | No | Moderate | Yes |
| 2 | 30 | 22.0 | 95 | No | Good | No |
| 3 | 35 | 25.4 | 105 | Yes | Poor | Yes |
| 4 | 60 | 33.5 | 155 | Yes | Poor | Yes |

```
1  df.head()
```

```
1  df.info()
```

```
1  # 1. Bar Chart: Jumlah pasien diabetes berdasarkan jenis diet
2  plt.figure(figsize=(8, 5))
3  sns.countplot(x='Diet', hue='Diabetes', data=df)
4  plt.title("Jumlah pasien diabetes berdasarkan jenis diet")
5  plt.xlabel("Diet Type")
6  plt.ylabel("Count")
7  plt.legend(title="Diabetes")
8  plt.show()
```

```python
# 2. Pie Chart: Persentase pasien dengan riwayat keluarga
family_history_counts = df['Family_History'].value_counts()
plt.figure(figsize=(5, 5))
family_history_counts.plot.pie(autopct='%1.1f%%', startangle=140, colors=['#ff9999','#66b3ff'])
plt.title("Persentase pasien dengan riwayat keluarga")
plt.ylabel("")
plt.show()
```

```python
# 3. Line Chart: Tren kadar gula darah berdasarkan usia
df_sorted = df.sort_values(by='Age')
plt.figure(figsize=(8, 5))
plt.plot(df_sorted['Age'], df_sorted['Blood_Sugar_Level'], marker='o')
plt.title("Tren kadar gula darah berdasarkan usia")
plt.xlabel("Age")
plt.ylabel("Blood Sugar Level (mg/dL)")
plt.grid(True)
plt.show()
```

```python
# 4. Scatter Plot: Hubungan antara BMI dan Kadar Gula Darah
plt.figure(figsize=(8, 5))
sns.scatterplot(x='BMI', y='Blood_Sugar_Level', hue='Diabetes', data=df)
plt.title("Hubungan antara BMI dan Kadar Gula Darah")
plt.xlabel("BMI")
plt.ylabel("Blood Sugar Level (mg/dL)")
plt.legend(title="Diabetes")
plt.show()
```

```python
# 5. Histogram: Distribusi BMI
plt.figure(figsize=(8, 5))
plt.hist(df['BMI'], bins=5, color='skyblue', edgecolor='black')
plt.title("Distribusi BMI")
plt.xlabel("BMI")
plt.ylabel("Frequency")
plt.show()
```

```python
1  # 6. Box Plot: Tingkat Gula Darah berdasarkan Riwayat Keluarga
2  plt.figure(figsize=(8, 5))
3  sns.boxplot(x='Family_History', y='Blood_Sugar_Level', data=df)
4  plt.title("Tingkat Gula Darah berdasarkan Riwayat Keluarga")
5  plt.xlabel("Family History")
6  plt.ylabel("Blood Sugar Level (mg/dL)")
7  plt.show()
```

```python
1  # 7. Heatmap: Korelasi antar variabel numerik
2  plt.figure(figsize=(8, 5))
3  sns.heatmap(df[['Age', 'BMI', 'Blood_Sugar_Level']].corr(), annot=True, cmap='coolwarm')
4  plt.title("Correlation Heatmap")
5  plt.show()
```

```python
1  # 8. Area Chart: Kasus diabetes kumulatif berdasarkan usia (untuk tujuan ilustrasi)
2  age_cum_cases = df[df['Diabetes'] == 'Yes'].groupby('Age').size().cumsum()
3  plt.figure(figsize=(8, 5))
4  plt.fill_between(age_cum_cases.index, age_cum_cases.values, color='lightgreen', alpha=0.6)
5  plt.title("Kasus Diabetes Kumulatif Berdasarkan Usia")
6  plt.xlabel("Age")
7  plt.ylabel("Cumulative Cases")
8  plt.grid(True)
9  plt.show()
```

```python
1  # Membuat DataFrame
2  df_encoded = pd.DataFrame(df)
```

```python
1  # Encode kolom categorical
2  df_encoded['Family_History'] = df_encoded['Family_History'].map({'Yes': 1, 'No': 0})
3  df_encoded['Diet'] = df_encoded['Diet'].map({'Poor': 0, 'Moderate': 1, 'Good': 2})
4  df_encoded['Diabetes'] = df_encoded['Diabetes'].map({'Yes': 1, 'No': 0})
```

```
1  df_encoded.head()
```

|   | Age | BMI | Blood_Sugar_Level | Family_History | Diet | Diabetes |
|---|-----|-----|-------------------|----------------|------|----------|
| 0 | 45  | 28.7 | 135 | 1 | 0 | 1 |
| 1 | 50  | 31.2 | 145 | 0 | 1 | 1 |
| 2 | 30  | 22.0 | 95  | 0 | 2 | 0 |
| 3 | 35  | 25.4 | 105 | 1 | 0 | 1 |
| 4 | 60  | 33.5 | 155 | 1 | 0 | 1 |

```
1  # mempersiapkan data untuk Categorical Naive Bayes
2  categorical_df = df_encoded.copy()
3  categorical_df['Age'] = pd.cut(categorical_df['Age'], bins=[0, 30, 40, 50, 60, 80], labels=[0, 1, 2, 3, 4])
4  categorical_df['BMI'] = pd.cut(categorical_df['BMI'], bins=[0, 18.5, 24.9, 29.9, 50], labels=[0, 1, 2, 3])
5  categorical_df['Blood_Sugar_Level'] = pd.cut(categorical_df['Blood_Sugar_Level'], bins=[0, 100, 125, 150, 200], labels=[0, 1, 2, 3])
```

```
1  categorical_df.head()
```

|   | Age | BMI | Blood_Sugar_Level | Family_History | Diet | Diabetes |
|---|-----|-----|-------------------|----------------|------|----------|
| 0 | 2 | 2 | 2 | 1 | 0 | 1 |
| 1 | 2 | 3 | 2 | 0 | 1 | 1 |
| 2 | 0 | 1 | 0 | 0 | 2 | 0 |
| 3 | 1 | 2 | 1 | 1 | 0 | 1 |
| 4 | 3 | 3 | 3 | 1 | 0 | 1 |

```python
1  # tentukan features dan target untuk CategoricalNB
2  X_categorical = categorical_df[['Age', 'BMI', 'Blood_Sugar_Level', 'Family_History', 'Diet']]
3  y_categorical = categorical_df['Diabetes']
```

```python
1  # membagi (split) dataset
2  X_train_cat, X_test_cat, y_train_cat, y_test_cat = train_test_split(X_categorical, y_categorical, test_size=0.3, random_state=42)
```

```python
1  # Creating and training the Categorical Naive Bayes model
2  cat_nb_model = CategoricalNB()
3  cat_nb_model.fit(X_train_cat, y_train_cat)
```

```python
1  # Predicting on the test set
2  y_pred_cat = cat_nb_model.predict(X_test_cat)
```

## 3. Simulasi Algoritma Naïve Bayes

```python
1  # Simulating with new input data (categorical values)
2  # Example: Age 45, BMI 27.5, Blood Sugar Level 125, Family History Yes, Diet Poor
3  # Convert these values to categorical codes
4  new_data_categorical = np.array([[3, 2, 1, 1, 0]])
5  new_prediction_cat = cat_nb_model.predict(new_data_categorical)
```

```python
1  (new_prediction_cat)
```

```python
1  if (new_prediction_cat[0]==0):
2      print ('Pasien tidak terkena diabetes')
3  else:
4      print ('Pasien terkena diabetes')
```

**4. Evaluasi Algoritma Naïve Bayes**

```python
1  # Calculating accuracy for Categorical Naive Bayes
2  cat_accuracy = accuracy_score(y_test_cat, y_pred_cat)
```

```python
1  # Generating the confusion matrix
2  cat_conf_matrix = confusion_matrix(y_test_cat, y_pred_cat)
```

```python
1  (cat_accuracy, cat_conf_matrix)
```

```python
1  # Menyimpan model menggunakan pickle
2  import pickle
3
4  filename = 'cat_diabetes_mod.pkl'
5  pickle.dump(cat_nb_model, open(filename, 'wb'))
```

**5. Aplikasi Menggunakan algoritma naïve bayes**

```python
1   import streamlit as st
2   import pickle
3   import numpy as np
4   import matplotlib.pyplot as plt
5
6   # Load the model from the .pkl file
7   with open('cat_diabetes_mod.pkl', 'rb') as model_file:
8       model = pickle.load(model_file)
9
10  # Streamlit app title
11  st.title("Categorical Naive Bayes - Diabetes Prediction")
12
13  # Sidebar inputs for user data
14  st.sidebar.header("Input Data Baru")
15  age = st.sidebar.selectbox("Age", ["<30", "30-40", "40-50", "50-60", "60+"])
16  bmi = st.sidebar.selectbox("BMI", ["<18.5", "18.5-24.9", "25-29.9", "30+"])
17  blood_sugar = st.sidebar.selectbox("Blood Sugar Level", ["<100", "100-125", "125-150", "150+"])
18  family_history = st.sidebar.selectbox("Family History of Diabetes", ["No", "Yes"])
19  diet = st.sidebar.selectbox("Diet", ["Good", "Moderate", "Poor"])
20
21  # Map inputs to numerical values (based on categories used for CategoricalNB)
22  age_map = {"<30": 0, "30-40": 1, "40-50": 2, "50-60": 3, "60+": 4}
23  bmi_map = {"<18.5": 0, "18.5-24.9": 1, "25-29.9": 2, "30+": 3}
24  blood_sugar_map = {"<100": 0, "100-125": 1, "125-150": 2, "150+": 3}
25  family_history_map = {"No": 0, "Yes": 1}
26  diet_map = {"Good": 2, "Moderate": 1, "Poor": 0}
27
28  # Convert user input to model input format
29  input_data = np.array([[age_map[age], bmi_map[bmi], blood_sugar_map[blood_sugar],
30                          family_history_map[family_history], diet_map[diet]]])
31
32  # Make prediction
33  prediction = model.predict(input_data)
34  prediction_prob = model.predict_proba(input_data)
35
36  # Display results
37  st.subheader("Prediction Result")
38  result = "Diabetes (Yes)" if prediction[0] == 1 else "No Diabetes"
39  st.write(f"The model predicts: **{result}**")
40
41  # Plot prediction probabilities
42  fig, ax = plt.subplots()
43  labels = ['No Diabetes', 'Diabetes']
44  ax.bar(labels, prediction_prob[0], color=['green', 'red'])
45  ax.set_ylabel('Probability')
46  ax.set_title('Prediction Probability')
47  st.pyplot(fig)
48
49  # Display input summary
50  st.subheader("Input Summary")
51  st.write(f"**Age**: {age}")
52  st.write(f"**BMI**: {bmi}")
53  st.write(f"**Blood Sugar Level**: {blood_sugar}")
54  st.write(f"**Family History**: {family_history}")
55  st.write(f"**Diet**: {diet}")
```