# Pokok Bahasan VII Naïve Bayes

Kode Pokok Bahasan: TIK.RPL03.005.00.01

## Deskripsi Pokok Bahasan:

Membahas bagaimana penerapan Algoritma Naïve Bayes pada dataset titanic.

No	Elemen Kompetensi	Indikator Kinerja	Jml Jam	Hal
1.	Menampilkan peluang dari kasus yang diberikan.	Mampu melakukan analisis terhadap peluang atas kejadian yang ditentukan.	1	12
2.	Menggunakan fungsi naiveBayes dalam memprediksi data	Mengimplementasikan fungsi naive bayes pada prediksi data		

#### **TUGAS PENDAHULUAN**

Hal yang harus dilakukan dan acuan yang harus dibaca sebelum praktikum:

- 1. Menginstal R pada PC masing-masing praktikan.
- 2. Menginstal R Studio pada PC masing-masing praktikan.

#### **DAFTAR PERTANYAAN**

1. Apa itu algoritma Naïve Bayes?

Algoritma Naive Bayes adalah algoritma yang mempelajari probabilitas suatu objek dengan ciri-ciri tertentu yang termasuk dalam kelompok/kelas tertentu. Singkatnya, ini adalah pengklasifikasi probabilistik.

2. Apa kegunaan Naïve Bayes?

Prediksi multi-kelas: Algoritma klasifikasi Naive Bayes dapat digunakan untuk memprediksi probabilitas posterior dari beberapa kelas variabel target. Klasifikasi teks: Karena fitur prediksi multi-kelas, Naive Bayes algoritma klasifikasi sangat cocok untuk klasifikasi teks.

3. Sebutkan tahapan dari proses algoritma Naïve Bayes!

Menghitung jumlah kelas/label. Menghitung jumlah kasus perkelas. Mengalikan semua hasil variable kelas. Membandingkan hasil perkelas.

## **TEORI SINGKAT**

Algoritma Naive Bayes merupakan sebuah metoda klasifikasi menggunakan metode probabilitas dan statistik yg dikemukakan oleh ilmuwan Inggris Thomas Bayes. Algoritma Naive Bayes 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.

Naive Bayes Classifier bekerja sangat baik dibanding dengan model classifier lainnya. Hal ini dibuktikan pada jurnal Xhemali, Daniela, Chris J. Hinde, and Roger G. Stone. "Naive Bayes vs. decision trees vs. neural networks in the classification of training web pages." (2009), mengatakan bahwa "Naïve Bayes Classifier memiliki tingkat akurasi yg lebih baik dibanding model classifier lainnya".

Keuntungan penggunan adalah bahwa metoda ini hanya membutuhkan jumlah data pelatihan (training data) yang kecil untuk menentukan estimasi parameter yg diperlukan dalam proses pengklasifikasian. Karena yg diasumsikan sebagai variabel independent, maka hanya varians dari suatu variabel dalam sebuah kelas yang dibutuhkan untuk menentukan klasifikasi, bukan keseluruhan dari matriks kovarians.

#### LAB SETUP

Hal yang harus disiapkan dan dilakukan oleh praktikan untuk menjalankan praktikum modul ini.

- 1. Menginstall library yang dibutuhkan untuk mengerjakan modul.
- 2. Menjalankan R Studio.

#### **ELEMEN KOMPETENSI I**

#### Deskripsi:

Menampilkan peluang dari kasus yang diberikan.

#### Kompetensi Dasar:

Mampu melakukan analisis terhadap peluang atas kejadian yang ditentukan.

#### Latihan

#### Penjelasan Singkat:

Pada latihan ini anda akan diminta untuk menampilkan summarize dari data menggunakan library yang disediakan oleh R.

#### Langkah-Langkah Praktikum:

1. Instal dan panggil package berikut jika belum terinstall

install.packages("tidyverse")
library("tidyverse") #for data wrangling tools
install.packages("titanic")
library("titanic")

2. Inisialisasi titanic\_train ke dalam variabel tdf.

tdf <- titanic\_train #training set of Titanic data head(tdf)

#### 3. Lakukan perintah berikut:

Compute the probability that a randomly selected passenger on the Titanic was female given that the passenger was at least 35 years old.

tdf %>%

summarize(prob = sum(Age >= 35 & Sex == "female", na.rm = TRUE)/sum(Age >= 35, na.rm = TRUE))

#### Output:

```
head(tdf)
  PassengerId Survived Pclass
                                                                                          Name
                                                                   Braund, Mr. Owen Harris
                                                                                                  male 22
                                                                                                                        0
                        0
                                1 Cumings, Mrs. John Bradley (Florence Briggs Thayer)
                                                                                                                        0
                                          Heikkinen, Miss. Laina
Futrelle, Mrs. Jacques Heath (Lily May Peel)
Allen, Mr. William Henry
                                                                                               female
                                                                                                                 0
                                                                                                                        0
                                                                                                         35
                                                                                                                        0
                                                                                               female
6
                                                                           Moran, Mr. James
              Ticket
                         Fare Cabin Embarked
           PC 17599 71.2833
3101282 7.9250
                                 C85
3 STON/02. 3101282
              113803 53.1000
                                C123
              373450
                      8.0500
              330877
6
                      8.4583
    summarize(prob = sum(Age >= 35 & Sex == "female", na.rm = TRUE)/sum(Age >= 35, na.rm = TRUE))
1 0.3446809
```

#### **ELEMEN KOMPETENSI II**

#### Deskripsi:

Menggunakan fungsi naïve bayes dalam memprediksi data.

## Kompetensi Dasar:

Mengimplementasikan fungsi naive bayes pada prediksi data

#### Latihan 1.2.1

#### Penjelasan Singkat:

Pada latihan ini anda akan diminta untuk mengimplementasikan naïve bayes pada kasus yang diberikan.

#### Langkah-Langkah Praktikum:

Gunakan <u>titanic.csv</u> yang berisi data 887 penumpang Titanic passengers. Kolom data menggambarkan survived (S), age (A), passenger-class (C), sex (G) and the fare paid (X). Hitung peluang bersyarat (conditional probability) di bawah ini

```
P(S= true | G=female)
```

```
P(S= \text{true} \mid G=\text{male})

P(S= \text{true} \mid C=1)

P(S= \text{true} \mid C=2)

P(S= \text{true} \mid C=3)

P(S= \text{true} \mid G=\text{female}, C=1) =

P(S= \text{true} \mid G=\text{female}, C=2) =

P(S= \text{true} \mid G=\text{male}, C=1) =

P(S= \text{true} \mid G=\text{male}, C=2) =

P(S= \text{true} \mid G=\text{male}, C=3) =

P(S=\text{true} \mid G=\text{male}, C=3) =
```

```
> #S=True/G=Female
> tdf %>%
    Summarize(prob = sum(Survived == "1" & Sex == "female", na.rm = TRUE)/sum(Survived == "1", na.rm = TRUE))
1 0 6812865
> #S=True/G=Male
> tdf %>%
    summarize(prob = sum(Survived == "1" & Sex == "male", na.rm = TRUE)/sum(Survived == "1", na.rm = TRUE))
        prob
1 0.3187135
> #S=True/C=1
> tdf %>%
    summarize(prob = sum(Survived == "1" & Pclass == "1", na.rm = TRUE)/sum(Survived == "1", na.rm = TRUE))
prob
1 0.3976608
   #S=True/C=2
> tdf %-%
+ summarize(prob = sum(Survived == "1" & Pclass == "2", na.rm = TRUE)/sum(Survived == "1", na.rm = TRUE))
prob
1 0.254386
  #S=True/C=3
tdf %%
summarize(prob = sum(Survived == "1" & Pclass == "3", na.rm = TRUE)/sum(Survived == "1", na.rm = TRUE))
prob
1 0.3479532
> #S=True/G=Female/C=1
> tdf %>%
     summarize(prob = sum(Survived == "1" & Sex == "female" & Pclass == "1", na.rm = TRUE)/sum(Survived == "1", na.rm = TRU
        prob
1 0.2660819
> #S=True/G=Female/C=2
> tdf %>%
+ summarize(prob = sum(Survived == "1" & Sex == "female" & Pclass == "2", na.rm = TRUE)/sum(Survived == "1", na.rm = TRUE))
1 0.2046784
> #S=True/G=Female/C=3
> tdf %>%
+ summarize(prob = sum(Survived == "1" & Sex == "female" & Pclass == "3", na.rm = TRUE)/sum(Survived == "1", na.rm = TRUE))
prob
1 0.2105263
> #S=True/G=male/C=1
> tdf %>%
    summarize(prob = sum(Survived == "1" & Sex == "male" & Pclass == "1", na.rm = TRUE)/sum(Survived == "1", na.rm = TRUE))
1 0.1315789
> #S=True/G=male/C=2
> tdf %>%
   Summarize(prob = sum(Survived == "1" & Sex == "male" & Pclass == "2", na.rm = TRUE)/sum(Survived == "1", na.rm = TRUE))
prob
1 0.0497076
> #S=True/G=male/C=3
> tdf %>%
    summarize(prob = sum(Survived == "1" & Sex == "male" & Pclass == "3", na.rm = TRUE)/sum(Survived == "1", na.rm = TRUE))
1 0.1374269
```

#### Jalankan perintah R di bawah ini :

```
# https://www.kaggle.com/brirush/naive-bayes-for-titanic
library(e1071)
train <- read.csv("F:/Kuliah Data Mining gasal 1819/Kaggle/Titanic/train.csv")
```

```
test <- read.csv("F:/Kuliah Data Mining gasal 1819/Kaggle/Titanic/test.csv")

BayesTitanicModel<-naiveBayes(as.factor(Survived)~., train)

BayesPrediction<-predict(BayesTitanicModel, test)
summary(BayesPrediction)
output<-data.frame(test$PassengerId, BayesPrediction)
str(output)
colnames(output)<-cbind("PassengerId", "Survived")
write.csv(output, file = 'Rushton_Solution.csv', row.names = F)
```

#### Output:

```
| https://www.kaggle.com/brirush/naive-bayes-for-titanic | library(e1071) | library(e1071)
```

## Berikan penjelasan terhadap output di atas

Mengeluarkan data dari File, mengecek data dari penumpang dan status nya survived atau tidak.

## Tugas: Kasus "playing golf"

## Data excelNaive

			Humadit		PlayGol
id	Outlook	Temp	у	Wndy	f
1	Rainy	Hot	High	FALSE	No
2	Rainy	Hot	High	TRUE	No
	Overcas				
3	t	Hot	High	FALSE	Yes
4	Sunny	Mild	High	FALSE	Yes
5	Sunny	Cool	Normal	FALSE	Yes
6	Sunny	Cool	Normal	TRUE	No
	Overcas				
7	t	Cool	Normal	TRUE	Yes
8	Rainy	Mild	High	FALSE	No
9	Rainy	Cool	Normal	FALSE	Yes
10	Sunny	Mild	Normal	FALSE	Yes
11	Rainy	Mild	Normal	TRUE	Yes
	Overcas				
12	t	Mild	High	TRUE	Yes
	Overcas				
13	t	Hot	Normal	FALSE	Yes
14	Sunny	Mild	High	TRUE	No

## Data excelNaiveTest

				Humadit	
id		Outlook	Temp	у	Wndy
	16	Rainy	Mild	Normal	TRUE

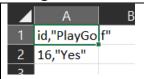
## Script R:

- > excelNaive <- read.delim('clipboard')
- > excelNaiveTest <- read.delim('clipboard')</pre>
- > excelNaiveModel <- naiveBayes(as.factor(PlayGolf)~., excelNaive)
- > excelNaivePredic <- predict(excelNaiveModel, excelNaiveTest)
- > summary(excelNaivePredic)
- > excelNaiveOutput <- data.frame(excelNaiveTest\$id, excelNaivePredic)
- > str(excelNaiveOutput)
- > colnames(excelNaiveOutput) <- cbind('id', 'PlayGolf')</pre>
- > write.csv(excelNaiveOutput, file = "rainy.csv", row.names = F)

### Output:

```
> #TUGAS
> excelNaive <- read.delim('clipboard')
> excelNaiveTest <- read.delim('clipboard')
> excelNaiveModel <- naiveBayes(as.factor(PlayGolf)~., excelNaive)
> excelNaivePredic <- predict(excelNaiveModel, excelNaiveTest)
> summary(excelNaivePredic)
No Yes
0 1
> excelNaiveOutput <- data.frame(excelNaiveTest$id, excelNaivePredic)
> str(excelNaiveOutput)
'data.frame': 1 obs. of 2 variables:
$ excelNaiveTest.id: int 16
$ excelNaiveTest.id: int 16
$ excelNaivePredic: Factor w/ 2 levels "No", "Yes": 2
> colnames(excelNaiveOutput) <- cbind('id', 'PlayGolf')
> write.csv(excelNaiveOutput, file = "rainy.csv", row.names = F)
> |
```

## Hitungan manual data test:



Setelah mengolah data dan survei terhadap data-data yang sudah ada, ketika id baru masuk dengan data yang ada maka terbukti lah data yang baru masuk seprti diatas

#### Sumber:

http://web.stanford.edu/class/archive/cs/cs109/cs109.1166/problem12.html

http://rstudio-pubs-static.s3.amazonaws.com/6595 b57093a21dfc46a4b3338cfee29ec61e.html

https://community.rstudio.com/t/conditional-probability-with-dplyr/5117

https://www.kaggle.com/brirush/naive-bayes-for-titanic

https://www.geeksforgeeks.org/naive-bayes-classifiers/

#### **CEK LIST**

Elemen	No Lotilean	Penyelesaian		
Kompetensi	No Latihan	Selesai	Tidak selesai	
1	1.1.1	✓		

2	1.2.1	✓	

# FORM UMPAN BALIK

Elemen Kompetensi	Tingkat Kesulitan	Tingkat Ketertarikan	Waktu Penyelesaian dalam menit
Memahami data pre- processing.	Sangat Mudah	Tidak Tertarik	
	Mudah	Cukup Tertarik	
	☐ ✓ Biasa	Tertarik	
	Sulit	✓ Sangat Tertarik	
	Sangat Sulit		
Mengimplementasika n pre-processing data.	Sangat Mudah	Tidak Tertarik	
	Mudah	Cukup Tertarik	
	☐ ✓ Biasa	Tertarik	
	Sulit	✓ Sangat Tertarik	
	Sangat Sulit		
Elemen Kompetensi	Tingkat Kesulitan	Tingkat Ketertarikan	Waktu Penyelesaian dalam menit
Menggunakan fungsi naiveBayes dalam	Sangat Mudah	Tidak Tertarik	
memprediksi data	Mudah	Cukup Tertarik	
	☐ ✓ Biasa	Tertarik	
	Sulit	✓ Sangat Tertarik	
	Sangat Sulit		