PENAMBANGAN DATA DAN INTELIGENSI BISNIS BEKERJA DENGAN MISSING VALUES



Mukhtada Billah Nasution F1E122037

PROGRAM STUDI SISTEM INFORMASI FAKULTAS SAINS DAN TEKNOLOGI UNIVERSITAS JAMBI

2024

I. Menyiapkan data dari BPS

Data diambil dari website BPS yaitu Rata-Rata Harga Unit Pembangunan Rumah oleh Perum Perumnas tahun 2017-2018. Berikut lampiran datanya:

•	1 [‡]	2017 [‡]	2018 [‡]
1	ACEH	121.60	115.72
2	SUMATERA UTARA	40.44	73.94
3	RIAU	123.01	NA
4	JAMBI	29.80	44.32
5	SUMATERA SELATAN	121.51	142.06
6	BENGKULU	30.67	79.59
7	LAMPUNG	54.21	136.19
8	KEP. RIAU	138.99	109.97
9	DKI JAKARTA	214.10	NA
10	JAWA BARAT	188.74	78.24
11	JAWA TENGAH	130.54	128.07
12	DI YOGYAKARTA	NA	NA
13	JAWA TIMUR	304.40	170.20
14	BANTEN	97.95	116.70
15	BALI	14.18	41.13
16	NUSA TENGGARA BARAT	38.37	56.27
17	KALIMANTAN BARAT	18.23	66.83
18	KALIMANTAN TENGAH	51.81	68.45
19	KALIMANTAN SELATAN	500.47	NA
20	KALIMANTAN UTARA	NA	NA
21	SULAWESI UTARA	83.76	108.05
22	SULAWESI TENGAH	39.18	NA
23	SULAWESI SELATAN	185.69	74.92
24	SULAWESI TENGGARA	86.75	NA
25	INDONESIA	158.37	93.52

Gambar: Data dari BPS.

II. Cek Missing Values

Cara cek missing values dilakukan dengan menggunakan kode:

is.na(data)

// data adalah variable yang menampung keseluruhan dataframe dari
data BPS

is.na adalah method yang berfungsi untuk *checking missing values* pada suatu *dataframe*. Method ini akan mengembalikan *array* dengan nilai masing-masing menjadi Boolean True atau False.

```
> is.na(data)
       2017
             2018
 [1,] FALSE FALSE
 [2,] FALSE FALSE
 [3,] FALSE TRUE
 [4,] FALSE FALSE
 [5,] FALSE FALSE
 [6,] FALSE FALSE
 [7,] FALSE FALSE
 [8,] FALSE FALSE
 [9,] FALSE TRUE
[10,] FALSE FALSE
[11,] FALSE FALSE
[12,] TRUE TRUE
[13,] FALSE FALSE
[14,] FALSE FALSE
[15,] FALSE FALSE
[16,] FALSE FALSE
[17,] FALSE FALSE
[18,] FALSE FALSE
[19,] FALSE
            TRUE
[20,]
     TRUE
             TRUE
[21,] FALSE FALSE
[22,] FALSE TRUE
[23,] FALSE FALSE
[24,] FALSE TRUE
[25,] FALSE FALSE
```

Visualisasi missing values

III. Cek jumlah Missing Values

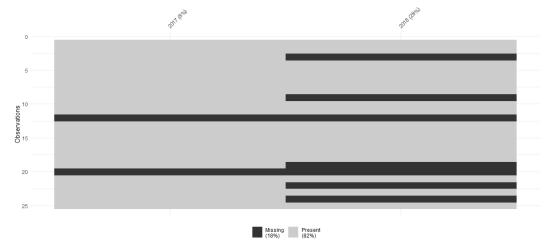
```
> sum(is.na(data))
[1] 9
```

Missing Value berjumlah 9 (Sembilan) buah. Sum adalah method yang berfungsi mencari akumulasi nilai tertentu, dalam hal ini akan mencari banyak nilai NA pada data.

IV. Visualisasi Missing Values

Visualisasi Missing Values dilakukan dengan menggunakn kode:

> visdat::vis miss(data)



Visdat::vis_miss() adalah suatu paket yang dapat memvisualisasikan missing data pada suatu dataframe. **Data** dijadikan sebagai argument pada method vis_miss() sehingga akan menampilkan graf keseluruhan data.

V. Hapus Missing Values

```
na.omit(data)
   A tibble: 18 × 2
2017 2018
       <db7>
                   <db7>
                  116.
       122.
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
        40.4
                    73.9
        29.8
                    44.3
                  142.
79.6
         30.7
         54.2
                  136.
       139.
                  110.
                    78.2
       189.
       131.
                  128.
       304.
                  170.
        98.0
                  117.
                   41.1
56.3
66.8
        14.2
         38.4
        18.2
         51.8
                    68.4
        83.8
                  108.
                    74.9
93.5
       186.
158.
```

na.omit adalah suatu method yang digunakan untuk mengembalikan dataframe yang bersih dari *missing values* (NA). Method ini mengembalikan dataframe.

```
> print(data, n = 25)
# A tibble: 25 x 2
    `2017` `2018`
                            > na.omit(data)
                            # A tibble: 18 \times 2
                                `2017` `2018`
    <db1>
           <db1>
    122.
           116.
                                 <db7>
                                          <db1>
     40.4
           73.9
                                 122.
                                          116.
    123.
            NA
                             2
                                  40.4
                                           73.9
     29.8
           44.3
                             3
                                  29.8
                                           44.3
    122.
           142.
            79.6
 6
     30.7
                            4
                                 122.
                                          142.
     54.2
           136.
                             5
                                           79.6
                                  30.7
 8
    139.
           110.
                            6
                                  54.2
                                          136.
9
    214.
            NA
                             7
                                 139.
                                          110.
    189.
            78.2
11
    131.
           128.
                             8
                                 189.
                                           78.2
     NA
           NA
                            9
                                 131.
                                          128.
    304.
           170.
13
                            10
                                 304.
                                          170.
     98.0 117.
14
                                  98.0
15
     14.2
           41.1
                            11
                                         117.
16
     38.4
            56.3
                            12
                                  14.2
                                           41.1
17
     18.2
            66.8
                           13
                                  38.4
                                           56.3
18
    51.8
            68.4
                            14
                                  18.2
                                           66.8
19
    500.
            NA
     NA
            NA
                            15
                                  51.8
                                           68.4
21
     83.8 108.
                            16
                                  83.8
                                          108.
22
     39.2
            NA
                            17
                                 186.
                                           74.9
23
    186.
            74.9
                            18
                                 158.
                                           93.5
24
     86.8
            NA
25
    158.
            93.5
                            >
```

Dataframe dengan (kiri) dan tanpa missing values (kanan)

VI. Mengganti Missing Values dengan mean

Pada bagian ini mengganti *Missing Values* (NA) dengan mean dari seluruh nilai pada kolom.

```
> print(data, n = 25)
 A tibble: 25 x 2
`2017` `2018`
    <db1>
            <db1>
    122.
            116.
     40.4
             73.9
    123.
             NA
     29.8
             44.3
    122.
            142.
     30.7
             79.6
     54.2 136.
 8
   139.
            110.
    214.
             NA
             78.2
10
    189.
11
    131.
            128.
12
            NA
     NA
            170.
13
    304.
14
     98.0 117.
15
     14.2
             41.1
16
     38.4
             56.3
17
     18.2
             66.8
18
     51.8
             68.4
19
    500.
             NA
     NA
             NA
21
     83.8 108.
     39.2
             NA
             74.9
    186.
24
     86.8
             NA
             93.5
    158.
```

Mean didapat berdasarkan dataframe yang telah dihilangkan missing valuesnya.

```
> m1 = mean(data$^2017^, na.rm=T)
> m2 = mean(data$^2018^, na.rm=T)
> m1
[1] 118.952
> m2
[1] 94.67611
> data$^2017^[is.na(data$^2017^)] = m1
> data$^2018^[is.na(data$^2018^)] = m2
```

data\$`2017`[is.na(data\$`2017`)] = m1 berarti data setiap index pada data pada kolom 2017 yang bernilai NA = True (berarti nilainya NA) akan diganti dengan m1, di mana m1 adalah 118.952. Begitu pula pada data\$`2018`[is.na(data\$`2018`)] = m2

Berikut adalah data yang telah disempurnakan:

```
> print(data, n=25)
# A tibble: 25 \times 2
   `2017` `2018`
     <db1>
             <db1>
    122.
            116.
 2
     40.4
              73.9
 3
    123.
              94.7
 4
              44.3
      29.8
 5
    122.
            142.
 6
      30.7
              79.6
      54.2
            136.
 8
    139.
            110.
 9
    214.
              94.7
10
    189.
              78.2
11
    131.
            128.
12
    119.
              94.7
13
    304.
            170.
14
      98.0
            117.
15
      14.2
              41.1
16
      38.4
              56.3
17
     18.2
              66.8
18
      51.8
              68.4
19
    500.
              94.7
20
    119.
              94.7
21
      83.8
            108.
22
              94.7
      39.2
23
    186.
              74.9
24
              94.7
      86.8
25
    158.
              93.5
```

Data yang missing valuenya telah diganti dengan mean masing-masing kolom.

Before After

```
> print(data, n = 25)
                               > print(data, n=25)
# A tibble: 25 × 2
`2017` `2018`
                                # A tibble: 25 x 2
`2017` `2018`
    <db1> <db1>
                                    <db1>
                                           <db1>
    122.
           116.
                                   122.
                                            116.
     40.4
             73.9
                                            73.9
                                    40.4
    123.
             NA
                                   123.
                                             94.7
     29.8
            44.3
                                4
                                     29.8
                                             44.3
    122.
           142.
                                5
                                   122.
                                            142.
    30.7
            79.6
                                6
                                     30.7
                                             79.6
     54.2 136.
                                     54.2
                                           136.
8
    139.
           110.
                                8 139.
                                           110.
    214.
            NA
                                9
                                   214.
                                             94.7
             78.2
    189.
                               10 189.
                                             78.2
11
    131.
           128.
                               11
                                   131.
                                           128.
12
            NA
                                             94.7
                               12
                                    119.
    304.
           170.
13
                               13
                                    304.
                                           170.
14
     98.0
           117.
                               14
                                     98.0
                                           117.
15
     14.2
             41.1
                               15
                                     14.2
                                             41.1
     38.4
16
             56.3
                               16
                                     38.4
                                             56.3
17
     18.2
             66.8
                               17
                                     18.2
                                             66.8
18
     51.8
             68.4
                               18
                                     51.8
                                             68.4
19
    500.
             NA
                               19
                                   500.
                                             94.7
     NA
             NA
                                             94.7
                                   119.
21
     83.8
           108.
                               21
                                     83.8
                                           108.
22
     39.2
             NA
                               22
                                     39.2
                                             94.7
             74.9
23
    186.
                                23
                                    186.
                                             74.9
24
     86.8
             NA
                                24
                                     86.8
                                             94.7
25
    158.
             93.5
                               25 158.
                                             93.5
```

VII. Pengisian missing values dengan median

```
summary(data)
       2017
                           2018
         : 14.18
                             : 41.13
 Min.
                    Min.
 1st Qu.: 39.81
Median : 97.95
                     1st Qu.: 69.82
                     Median : 86.56
         :120.56
                             : 94.68
 Mean
                     Mean
 3rd Qu.:148.68
                     3rd Qu.:116.45
         :500.47
                     Max.
NA's
                             :170.20
 Max.
 NA's
         : 2
> md1 = 97.95
 md2 = 86.56
```

Method summary adalah suatu method yang digunakan untuk mencari informasi detail dari suatu dataframe. Seperti yang dilihat method ini mengembalikan berbagai informasi seeprti min, max, 1st Quarter, Mean, 3rd Quarter, dan jumlah NA. Median setiap kolom dimasukkan ke variable md1 (2017) dan md2(2018)

Untuk mengganti *missing values* dengan median, caranya kurang lebih sama dengan cara mengganti *missing values* dengan *mean*.

```
data$ 2017 [is.na(data$ 2017)] = md1
data$ 2018 [is.na(data$ 2018)] = md2
  <db7>
              <db1>
     122.
              116.
               73.9
      40.4
     123.
               86.6
      29.8
               44.3
     122.
              142.
 6
7
               79.6
      30.7
      54.2
              136.
 8
     139.
              110.
 9
     214.
               86.6
               78.2
10
     189.
11
     131.
              128.
12
      98.0
               86.6
13
     304.
              170.
14
      98.0
              117.
15
      14.2
               41.1
16
      38.4
               56.3
      18.2
17
               66.8
18
      51.8
               68.4
19
20
21
22
23
24
25
     500.
               86.6
      98.0
               86.6
      83.8
              108.
               86.6
74.9
      39.2
     186.
      86.8
               86.6
     158.
               93.5
```

Seperti yang dapat dilihat setiap nilai NA telah diganti dengan median masingmasing kolom.

Before After

```
> print(data, n = 25)
                           # A tibble: 25 × 2
`2017` `2018`
# A tibble: 25 \times 2
    `2017`
            2018`
                                <db1>
                                       <db1>
    <db1>
            <db1>
                              122.
                                       116.
    122.
            116.
                                40.4
                                        73.9
    40.4
            73.9
                               123.
                                        86.6
    123.
             NA
                                        44.3
             44.3
                            4
                                 29.8
 4
     29.8
                            5
                               122.
                                       142.
    122.
            142.
                            6
                                 30.7
                                        79.6
     30.7
 6
            79.6
                            7
                                 54.2
                                       136.
     54.2
           136.
                            8
                               139.
                                       110.
 8
    139.
            110.
                                        86.6
                            9
                               214.
9
    214.
             NA
                           10 189.
                                        78.2
10
    189.
             78.2
                                       128.
                           11 131.
11
    131.
            128.
                           12
                                 98.0
                                        86.6
12
     NA
            NA
                           13
                                304.
                                       170.
13
    304.
            170.
                           14
                                 98.0
                                       117.
     98.0
14
           117.
                           15
                                 14.2
                                        41.1
15
     14.2
             41.1
                           16
                                 38.4
                                        56.3
     38.4
             56.3
16
                           17
                                 18.2
                                        66.8
17
     18.2
             66.8
                                        68.4
                           18
                                 51.8
18
     51.8
             68.4
                           19
                                500.
                                        86.6
19
    500.
             NΑ
                           20
                                 98.0
                                        86.6
     NA
             NA
                           21
                                 83.8
                                       108.
21
     83.8
           108.
                           22
                                        86.6
                                 39.2
22
     39.2
             NA
                           23
                                186.
                                        74.9
23
    186.
             74.9
                           24
                                 86.8
                                        86.6
24
     86.8
             NA
                           25
                                158.
                                        93.5
25
    158.
             93.5
                           >
< I
```