

Computational Document

체질량지수(BMI) - R마크다운(.rmd)

true

2021-08-02

Contents

1	비즈니스 설명	1
2	데이터	2
2.1	데이터 사전	2
2.2	데이터 가져오기	2
3	탐색적 데이터 분석	3
3.1	요약 통계	3
3.2	시각화	3
4	예측모형 - BMI 예측	4
5	BMI 예측	5

1 비즈니스 설명

캐글, “500 Person Gender-Height-Weight-Body Mass Index - Height and Weight random generated, Body Mass Index Calculated”에서 데이터를 바탕으로 고객에게 체중과 키 정보만 제공하면 체질량 지수(Body Mass Index, BMI)를 예측하는 모형을 개발하여 고객이 궁금해하는 서비스를 개발하고자 한다.

체질량 지수(體質量指數, Body Mass Index, BMI)는 인간의 비만도를 나타내는 지수로, 체중과 키의 관계로 계산된다. 키가 t 미터, 체중이 w 킬로그램일 때, BMI는 다음이 수식으로 표현된다. (키의 단위가 센티미터가 아닌 미터임에 유의해야 한다.)

$$BMI = \frac{w}{t^2}$$

체질량지수 (BMI지수)로 과체중 혹은 비만을 판정하는 한국 사례 ¹

구분	BMI 지수
고도 비만	40 이상
중등도 비만 (2단계 비만)	35 - 39.9
경도 비만 (1단계 비만)	30 - 34.9
과체중	25 - 29.9
정상	18.5 - 24.9
저체중	18.5 미만

¹ 위키백과, “체질량 지수” (2019-04-22 접근함)

2 데이터

캐글, “500 Person Gender-Height-Weight-Body Mass Index - Height and Weight random generated, Body Mass Index Calculated”에서 데이터를 바탕으로 고객에게 체중과 키 및 라벨 데이터 **index**가 준비되어 있어 키와 몸무게를 통해 BMI 예측한다.

2.1 데이터 사전

- Gender : Male / Female
- Height : Number (cm)
- Weight : Number (Kg)
- Index :
 - 0 : Extremely Weak
 - 1 : Weak
 - 2 : Normal
 - 3 : Overweight
 - 4 : Obesity
 - 5 : Extreme Obesity

2.2 데이터 가져오기

캐글에서 내려받은 원본 데이터를 살펴본다.

```
library(tidyverse)
```

```
# bmi_dat <- read_csv("https://raw.githubusercontent.com/statklee/author_carpentry_kr/gh-pages/data/500")
```

```
bmi_dat <- read_csv("data/bmi_dat.csv")
```

```
glimpse(bmi_dat)
```

```
Rows: 500
```

```
Columns: 4
```

```
$ Gender <chr> "Male", "Male", "Female", "Female", "Male", "Male", "Male", "Male", "Male", "F~
```

```
$ Height <dbl> 174, 189, 185, 195, 149, 189, 147, 154, 174, 169, 195, 159, 192, 155, 191, 153~
```

```
$ Weight <dbl> 96, 87, 110, 104, 61, 104, 92, 111, 90, 103, 81, 80, 101, 51, 79, 107, 110, 12~
```

```
$ Index <dbl> 4, 2, 4, 3, 3, 3, 5, 5, 3, 4, 2, 4, 3, 2, 2, 5, 5, 5, 5, 5, 5, 5, 4, 5, 2, 3, ~
```

```
bmi_dat %>%
```

```
DT::datatable()
```

Show entries

Search:

	Gender	Height	Weight	Index
1	Male	174	96	4
2	Male	189	87	2
3	Female	185	110	4
4	Female	195	104	3
5	Male	149	61	3
6	Male	189	104	3
7	Male	147	92	5
8	Male	154	111	5
9	Male	174	90	3
10	Female	169	103	4

Showing 1 to 10 of 500 entries

Previous 2 3 4 5 ... 50 Next

3 탐색적 데이터 분석

3.1 요약 통계

```
bmi_df <- bmi_dat %>%
  mutate(Index = factor(Index, levels = c(0,1,2,3,4,5), labels = c(" ", " ", " ", " ", " ", " ")),
         Gender = factor(Gender, levels = c("Male", "Female")))

bmi_df %>%
  group_by(Index) %>%
  summarise(
    = mean(Height),
    = mean(Weight))
```

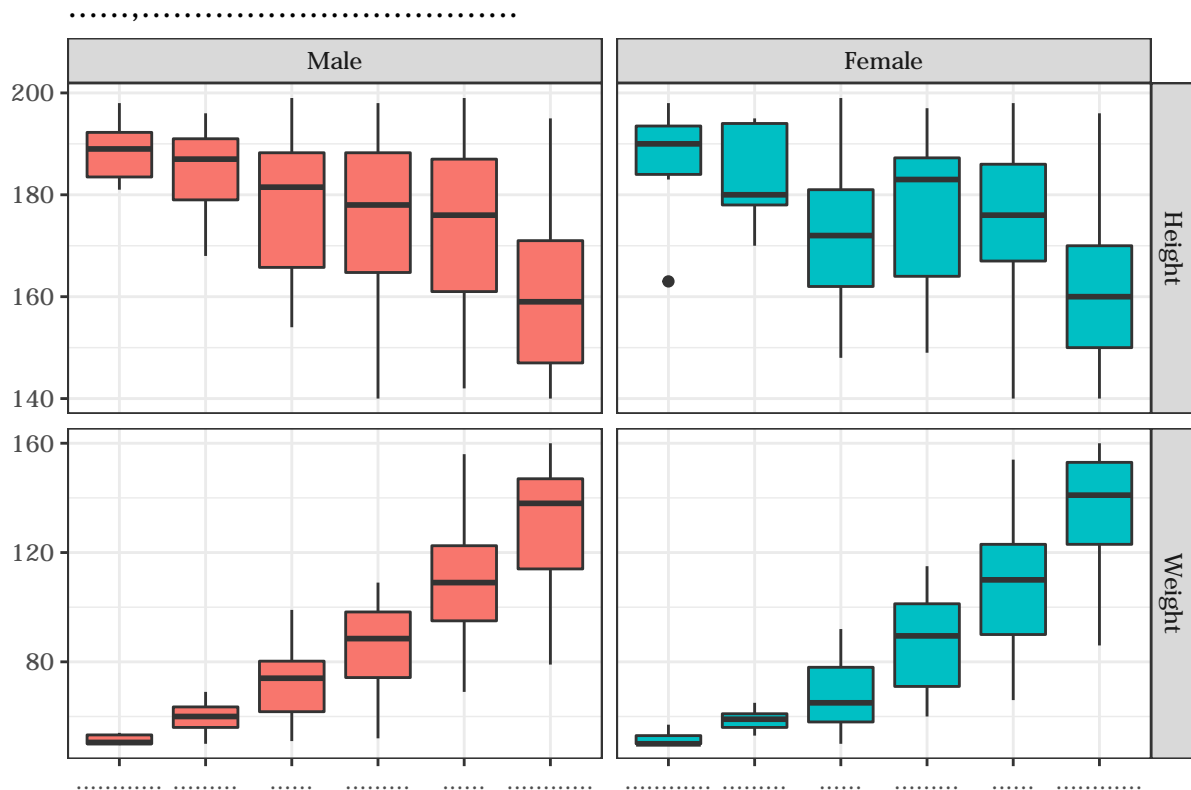
Index	평균키	평균체중
극저체중	187.5385	51.69231
저체중	184.7727	59.40909
정상	174.2609	69.08696
과체중	175.9853	86.88235
비만	173.8769	107.95385
고도비만	160.9798	132.88889

3.2 시각화

```
library(extrafont)
loadfonts()

bmi_df %>%
  gather( , , -Gender, -Index) %>%
  ggplot(aes(x=Index, y= , fill=Gender)) +
    geom_boxplot(show.legend = FALSE) +
    facet_grid( ~ Gender, scales="free") +
```

```
labs(x="", y="",
      title=" ",
      theme_bw(base_family = "NanumGothic"))
```



4 예측모형 - BMI 예측

$$\text{BMI 그룹} = f(\text{성별, 키, 몸무게}) + \epsilon$$

BMI 그룹: “극저체중”, “저체중”, “정상”, “과체중”, “비만”, “고도비만”

```
# 0. -----
library(caret)
library(doSNOW)

set.seed(777)

# 1. -----
# bmi_df

# 2. -----

# 3. -----
## 3.1.
num_cores <- parallel::detectCores()
start_time <- Sys.time()

cl <- makeCluster(num_cores, type = "SOCK")
registerDoSNOW(cl)
```

```
## 3.2. vs /
train_test_index <- createDataPartition(bmi_df$Index, p = 0.7, list = FALSE)

train <- bmi_df[train_test_index, ]
test <- bmi_df[-train_test_index, ]

## 3.3. / -----
cv_folds <- createMultiFolds(train$Index, k = 10, times = 5)
cv_ctrl <- trainControl(method = "cv", number = 10,
                        index = cv_folds,
                        verboseIter = TRUE)

## 3.2.
### ranger
gc_ranger_model <- train(Index ~., train,
                        method = "ranger",
                        tuneLength = 7,
                        trControl = cv_ctrl)
```

note: only 2 unique complexity parameters in default grid. Truncating the grid to 2 .

Aggregating results

Selecting tuning parameters

Fitting mtry = 3, splitrule = extratrees, min.node.size = 1 on full training set

```
# 4. -----
gc_pred_class <- predict(gc_ranger_model, newdata = test, type="raw")
## -----
bmi_conf <- confusionMatrix(gc_pred_class, test$Index)

bmi_conf$table
```

	Reference					
Prediction	3	0	0	0	0	0
3	3	0	0	0	0	0
0	0	6	0	0	0	0
0	0	0	19	1	0	0
0	0	0	1	19	2	0
0	0	0	0	0	35	5
0	0	0	0	0	2	54

```
cat(" : ", scales::percent(bmi_conf$overall[["Accuracy"]]))
```

```
: 93%
```

```
stopCluster(cl)
```

5 BMI 예측

```
bmi_test_dat <- tribble(
  ~"Gender", ~"Height", ~"Weight",
  "Male", 149, 61,
  "Female", 172, 67
)
```

```
predict(gc_ranger_model, newdata = bmi_test_dat, type="raw")
```

```
[1]
```

```
Levels:
```