

Red and Black 170303 : id Masked

coop711

2018-03-26

Data

```
class_roll <- read.table("../data/class_roll_masked.txt",
                        header = TRUE,
                        stringsAsFactors = FALSE,
                        encoding = "UTF-8")

str(class_roll)
```

```
## 'data.frame':   160 obs. of  6 variables:
##  $ dept   : chr   "○○학과" "○○학과" "○○학과" "○○학과" ...
##  $ id     : int   20119999 20119999 20179999 20149999 20169999 20129999 20149999 20169999 20179999 2
0129999 ...
##  $ name   : chr   "강○○" "강○○" "강○○" "강○○" ...
##  $ year   : int   4 4 1 4 2 3 4 2 1 3 ...
##  $ email  : chr   "user_name@naver.com" "user_name@hanmail.net" "user_name@naver.com" "user_name@han
mail.net" ...
##  $ cell_no: chr   "010-9164-xxxx" "010-8574-xxxx" "010-6435-xxxx" "010-2066-xxxx" ...
```

Randomization

```
# set.seed(107)
N <- nrow(class_roll)
class_roll$group <- sample(1:N) %% 2 + 1
class_roll$group <- factor(class_roll$group,
                          labels = c("Red", "Black"))

red_id <- which(class_roll$group == "Red")
black_id <- which(class_roll$group == "Black")
```

학번

```
ID_16 <- factor(ifelse(substr(class_roll$id, 1, 4) >= 2016,
                          "younger_16", "older_16"),
                levels = c("younger_16", "older_16"))
kable(table("그룹" = class_roll$group,
            "16학번 기준" = ID_16))
```

	younger_16	older_16
Red	46	34
Black	41	39

```
ID_15 <- factor(ifelse(substr(class_roll$id, 1, 4) >= 2015,
                          "younger_15", "older_15"),
                levels = c("younger_15", "older_15"))
kable(table("그룹" = class_roll$group,
            "15학번 기준" = ID_15))
```

	younger_15	older_15
Red	53	27
Black	47	33

```
ID_14 <- factor(ifelse(substr(class_roll$id, 1, 4) >= 2014,
                          "younger_14", "older_14"),
                levels = c("younger_14", "older_14"))
kable(table("그룹" = class_roll$group,
            "14학번 기준" = ID_14))
```

	younger_14	older_14
Red	62	18
Black	58	22

```
ID_13 <- factor(ifelse(substr(class_roll$id, 1, 4) >= 2013,
                          "younger_13", "older_13"),
                levels = c("younger_13", "older_13"))
kable(table("그룹" = class_roll$group,
            "13학번 기준" = ID_13))
```

	younger_13	older_13
Red	75	5
Black	71	9

email 서비스업체

```
email_list <- strsplit(class_roll$email, "@", fixed = TRUE)
mail_com <- sapply(email_list, `[`, 2)
kable(table("그룹" = class_roll$group,
            "e-mail" = mail_com))
```

	daum.net	gmail.com	hanmail.net	nate.com	naver.com
Red	2	4	6	5	63
Black	0	2	3	2	72

성씨 분포

```
f_name <- substring(class_roll$name,
                    first = 1, last = 1)
kable(table("Group" = class_roll$group,
            "Family Name" = f_name))
```

	강	고	구	권	김	나	명	문	박	반	방	배	서	성	손	송	신	심	안	양	우	유	윤	이	임	장	전	정	조	차	최	하	한
Red	2	0	0	1	18	1	0	0	8	0	0	2	2	1	2	1	2	1	4	0	0	3	4	12	1	1	2	4	2	0	4	0	1
Black	4	1	1	3	18	0	1	1	5	1	1	0	2	0	0	2	1	0	1	1	1	2	1	9	2	2	0	4	5	2	5	1	1

많이 나오는 성씨

```
f_name_f <- factor(ifelse(f_name %in% c("김", "이", "박"),
                          f_name, "기타"),
                  levels = c("김", "이", "박", "기타"))
kable(table("Group" = class_roll$group,
            "Family Name" = f_name_f))
```

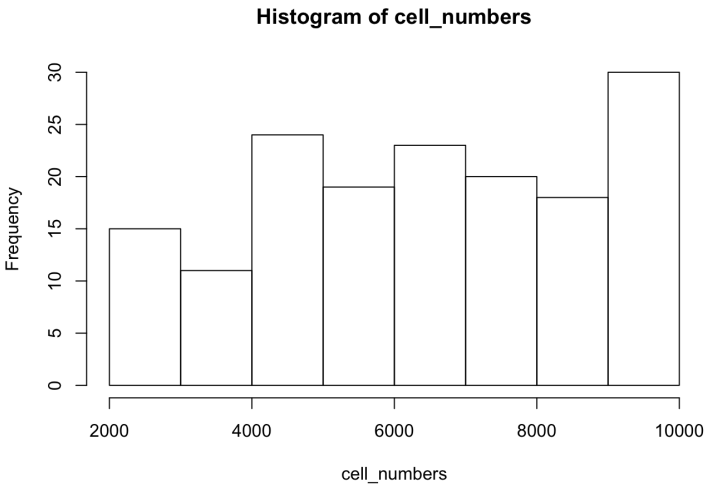
	김	이	박	기타
Red	18	12	8	42
Black	18	9	5	48

전화번호의 분포

```
cell_numbers <- sapply(substr(class_rol1$cell_no, 5, 8),
  as.numeric)
# cut_label <- c("1000~1999", "2000~2999", "3000~3999", "4000~4999", "5000~5999", "6000~6999",
#               "7000~7999", "8000~8999", "9000~9999")
cut_label <- paste(paste0(1:9, "000"), paste0(1:9, "999"), sep = "~")
kable(t(table(cut(cell_numbers,
  labels = cut_label,
  breaks = seq(1000, 10000, by = 1000))))))
```

1000~1999	2000~2999	3000~3999	4000~4999	5000~5999	6000~6999	7000~7999	8000~8999	9000~9999
0	15	11	24	19	23	20	18	30

```
hist(cell_numbers)
```



출석부에서 8명 비복원 랜덤 표집

```
# set.seed(1)
kable(class_rol1[sample(1:nrow(class_rol1), size = 8), ])
```

	dept	id	name	year	email	cell_no	group
145	○○학과	20179999	차○○	1	user_name@naver.com (mailto:user_name@naver.com)	010-8616-xxxx	Black
85	○○학과	20139999	안○○	3	user_name@naver.com (mailto:user_name@naver.com)	010-5030-xxxx	Red
119	○○학과	20169999	이○○	2	user_name@naver.com (mailto:user_name@naver.com)	010-4032-xxxx	Red
150	○○학과	20169999	최○○	2	user_name@naver.com (mailto:user_name@naver.com)	010-6379-xxxx	Red
63	○○학과	20179999	박○○	1	user_name@nate.com (mailto:user_name@nate.com)	010-6219-xxxx	Red
147	○○학과	20139999	최○○	3	user_name@naver.com (mailto:user_name@naver.com)	010-9079-xxxx	Black
22	○○학과	20139999	김○○	2	user_name@naver.com (mailto:user_name@naver.com)	010-7221-xxxx	Black
29	○○학과	20149999	김○○	3	user_name@daum.net (mailto:user_name@daum.net)	010-4039-xxxx	Red

set.seed() 의 용법

set.seed() 를 이용하면 랜덤넘버에 의존하는 실험을 재현할 수 있다. 다음 코드를 반복 수행하거나 다른 사람들의 수행결과와 비교해 보라.

세 결과가 모두 다른 경우

```
sample(1:6, size = 2)

## [1] 6 1

sample(1:6, size = 2)

## [1] 3 2

sample(1:6, size = 2)

## [1] 1 6
```

세 번의 수행 결과가 똑같이 반복되는 경우

```
set.seed(1)
sample(1:6, size = 2)

## [1] 2 6

sample(1:6, size = 2)

## [1] 4 5

sample(1:6, size = 2)

## [1] 2 5
```

동일한 결과를 반복적으로 얻는 경우

```
set.seed(1)
sample(1:6, size = 2)

## [1] 2 6

set.seed(1)
sample(1:6, size = 2)

## [1] 2 6

set.seed(1)
sample(1:6, size = 2)

## [1] 2 6
```