

# Red and Black 170303 : id Masked

MYUNG BIN KWAK

2020-04-04

## Data

```
class_rol1 <- read.table("../data/class_rol1_masked.txt",
                        header = TRUE,
                        stringsAsFactors = FALSE,
                        encoding = "CP949")

str(class_rol1)
```

```
## 'data.frame':    160 obs. of  6 variables:
## $ dept   : chr   "○○학과" "○○학과" "○○학과" "○○학과" ...
## $ id     : int  20119999 20119999 20179999 20149999 20169999 20129999 20149999 20169999 20179999 20129999 ...
## $ 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@hanmail.net" ...
## $ cell_no: chr   "010-1164-xxxx" "010-1174-xxxx" "010-1135-xxxx" "010-1166-xxxx" ...
```

## Randomization

```
set.seed(1)
N <- nrow(class_rol1)
class_rol1$group <- sample(1:N) %% 2
class_rol1$group <- factor(class_rol1$group,
                          labels = c("Red", "Black"))
red_id <- which(class_rol1$group == "Red")
black_id <- which(class_rol1$group == "Black")
```

## 학번

```
ID_16 <- factor(ifelse(substr(class_rol1$id, 1, 4) >= 2016,
                          "younger_16", "older_16"),
                levels = c("younger_16", "older_16"))
kable(table("그룹" = class_rol1$group,
            "16학번 기준" = ID_16))
```

	younger_16	older_16
Red	46	34
Black	41	39

```
ID_15 <- factor(ifelse(substr(class_rol1$id, 1, 4) >= 2015,
                          "younger_15", "older_15"),
                levels = c("younger_15", "older_15"))
kable(table("그룹" = class_rol1$group,
            "15학번 기준" = ID_15))
```

	younger_15	older_15
Red	54	26
Black	46	34

```
ID_14 <- factor(ifelse(substr(class_rol1$id, 1, 4) >= 2014,
                          "younger_14", "older_14"),
                levels = c("younger_14", "older_14"))
kable(table("그룹" = class_rol1$group,
            "14학번 기준" = ID_14))
```

	younger_14	older_14
Red	63	17
Black	57	23

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

	younger_13	older_13
Red	75	5
Black	71	9

## email 서비스업체

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

	daum.net	gmail.com	hanmail.net	nate.com	naver.com
Red	1	3	3	4	69
Black	1	3	6	3	66

## 성씨 분포

```
f_name <- substring(class_rolld$name,
                    first = 1, last = 1)
kable(table("Group" = class_rolld$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	황
Red	2	0	1	1	19	1	1	1	6	1	0	1	2	1	1	1	3	1	3	0	1	3	1	12	0	2	1	1	3	1	5	1	1	2
Black	4	1	0	3	17	0	0	0	7	0	1	1	2	0	1	2	0	0	2	1	0	2	4	10	2	1	1	7	4	1	4	0	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_rolld$group,
            "Family Name" = f_name_f))
```

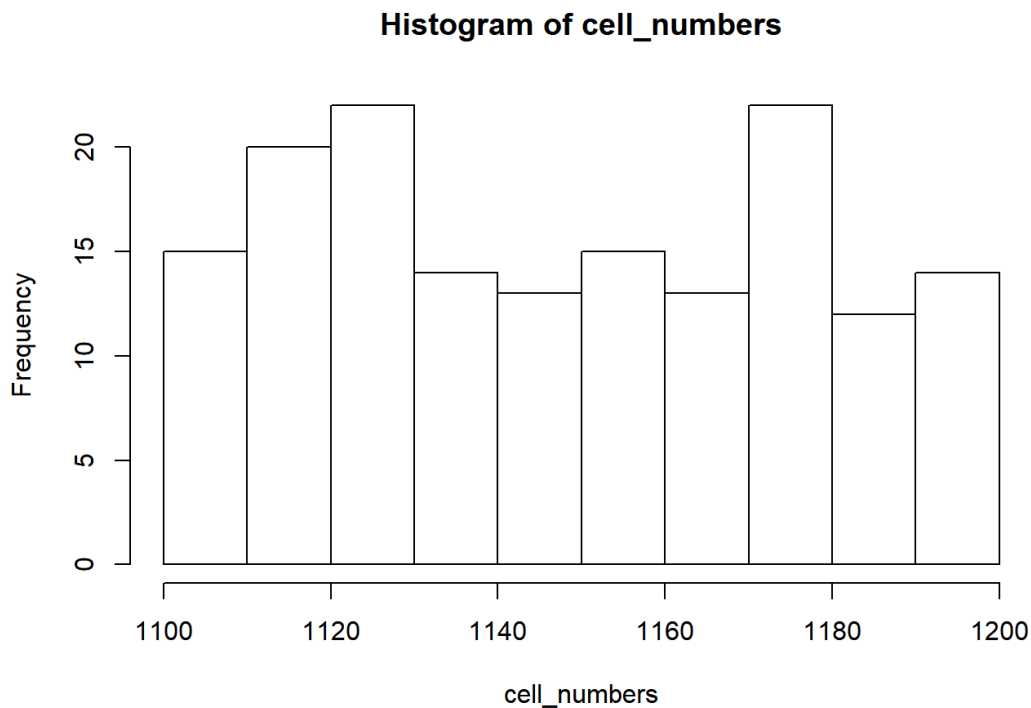
	김	이	박	기타
Red	19	12	6	43
Black	17	10	7	46

## 전화번호의 분포

```
cell_numbers <- sapply(substr(class_rol1$cell_no, 5, 8),
                      as.numeric)
# cut_label <- c("1100~1109", "1110~1119", "1120~1129", "1130~1139", "1140~1149", "1150~1159",
#               "1160~1169", "1170~1179", "1180~1189", "1190~1199")
cut_label <- paste(paste0(seq(1100, 1190, by = 10)), paste0(seq(1109, 1199, by = 10)), sep = "~")
kable(t(table(cut(cell_numbers,
                  labels = cut_label,
                  breaks = seq(1100, 1200, by = 10))))))
```

1100~1109	1110~1119	1120~1129	1130~1139	1140~1149	1150~1159	1160~1169	1170~1179	1180~1189	1190~1199
14	20	22	14	13	15	13	22	12	14

```
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
140	○○학과	20169999	조	○○	2	user_name@naver.com (mailto:user_name@naver.com)	010-1125-xxxx	Black
126	○○학과	20139999	장	○○	3	user_name@naver.com (mailto:user_name@naver.com)	010-1173-xxxx	Black
14	○○학과	20119999	김	○○	3	user_name@naver.com (mailto:user_name@naver.com)	010-1109-xxxx	Black
116	○○학과	20119999	이	○○	4	user_name@naver.com (mailto:user_name@naver.com)	010-1152-xxxx	Black
16	○○학과	20139999	김	○○	2	user_name@naver.com (mailto:user_name@naver.com)	010-1195-xxxx	Red
15	○○학과	20169999	김	○○	2	user_name@naver.com (mailto:user_name@naver.com)	010-1164-xxxx	Red
130	○○학과	20129999	정	○○	4	user_name@naver.com (mailto:user_name@naver.com)	010-1111-xxxx	Black
65	○○학과	20179999	반	○○	1	user_name@naver.com (mailto:user_name@naver.com)	010-1190-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] 3 1
```

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

```
## [1] 6 1
```

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

```
## [1] 4 5
```

세 번의 수행 결과가 똑같이 반복되는 경우

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

```
## [1] 1 4
```

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

```
## [1] 1 2
```

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

```
## [1] 5 3
```

동일한 결과를 반복적으로 얻는 경우

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

```
## [1] 1 4
```

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

```
## [1] 1 4
```

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

```
## [1] 1 4
```