# ABOUT DLOOKR PACKAGE

choonghyun ryu

# HISTORY

- who : choonghyun ryu

- when : after RStudio conference

- where : home

- what : R package

- why : motivation

- how : develop R package

# ABOUT

- Package: dlookr

- Type: Package

- Title: Tools for Data Exploration & Data Transformation

- Version: 0.3.0

- Authors@R: c(person("Choonghyun", "Ryu", email = "choonghyun.ryu@gmail.com", role = c("aut", "cre")))

- Description: A collection of tools that support data diagnosis, exploration, and transformation.

- License: GPL-2

# DATA DIAGNOSIS

# DATA EXPLORATION

# DATA TRANSFORMATION



**Transformation Information Report**

Report by dlookr package

2018-04-16

### 1.1.3 Urban

**Impute missing values with mode**

Table 1.7: Descriptive Statistics : Urban with 'mode'

|     | original | imputation | original_percent | imputation_percent |
|-----|----------|------------|------------------|--------------------|
| No  | 117      | 117        | 29.25            | 29.25              |
| Yes | 283      | 278        | 70.75            | 69.50              |
| NA  | 0        | 5          | 0.00             | 1.25               |

**Information of Imputation (before vs after)**



Figure 1.7: Urban - mode

"**pipe**를 이용한 응용,
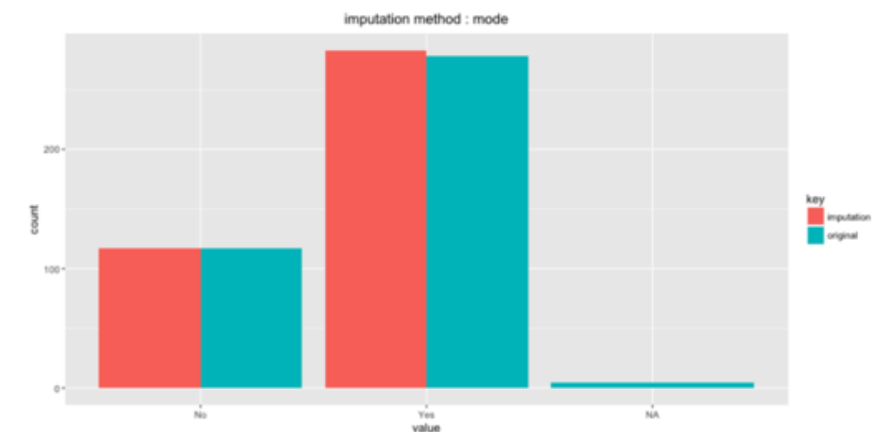**tidyverse**와의 궁합."

# BUCKET LIST

- 작은 소망

  hexbin 스티커를
  노트북에 붙이는 것