

# From poor data To training data

...

JeongMin Kwon



# JeongMin Kwon

Job Title:

Data Scientist

Interests:

Data Problem Solving, R, Python,  
translation,  
and many more.

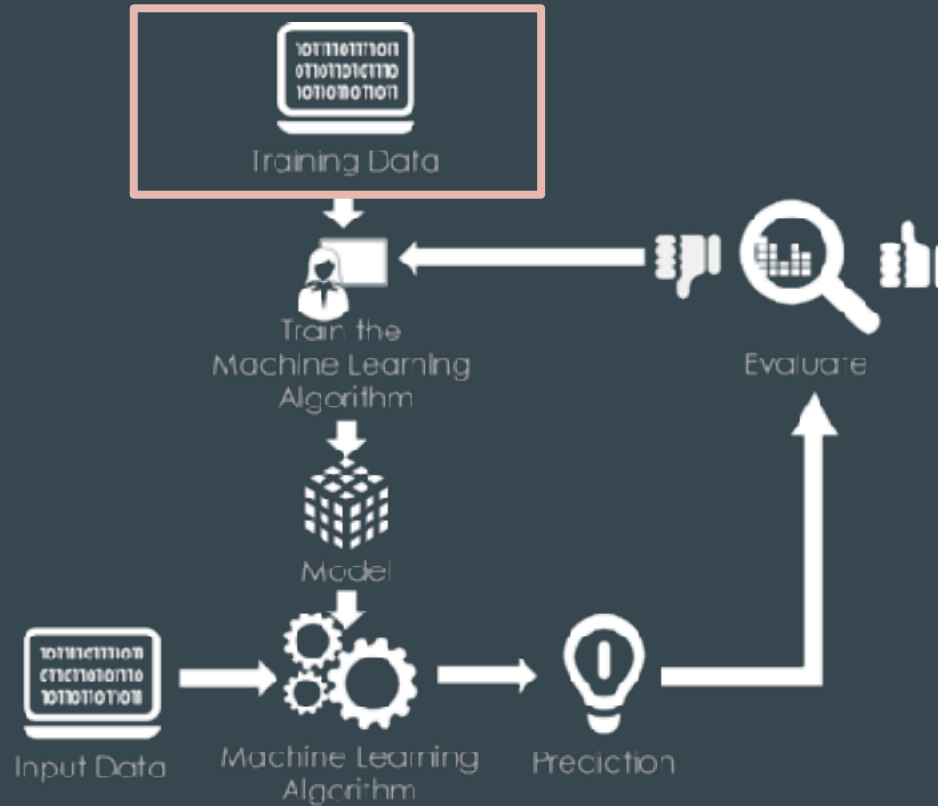
Contact:

cojette@gmail.com

/ @mazycat

---

# Machine learning services with Data



# Machine Learning Dataset Condition

- Clean
- Consistent
- Simple type
- Clear predictors and dependant variables
- Properly scaled and discretized

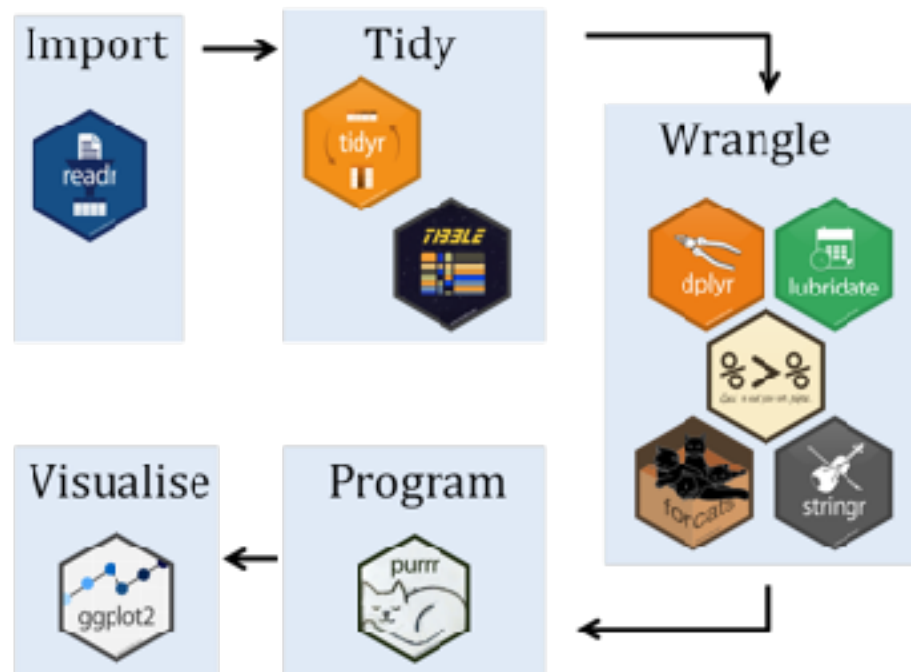
# There are lots of poor data...

**Tidy datasets are all alike  
but every messy dataset is messy in its own way**

- Raw and messy file (Raw mine for data)
- Data from various different sources and repositories
- Too large datasets
- Too various information or columns
- Different categories and criteria
- Meaningless text and unstructured data
- Lots of errata, omitted fields, various formats and missing values
- etc.

# Tidyverse

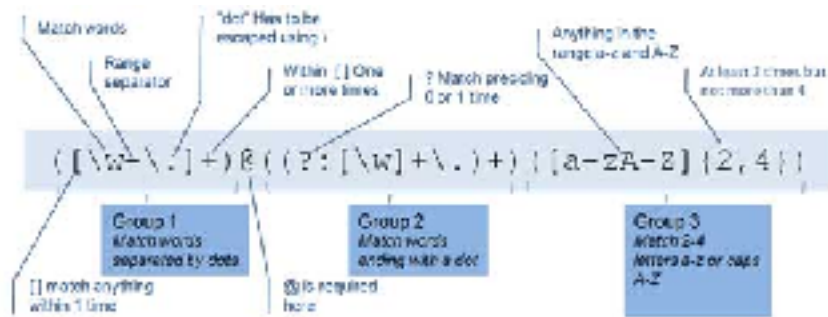
## Basic Data Wrangling



# Regular Expression

## Basic Data Wrangling

- Text Pattern control
- Text Cleansing
- Default: extended regex
  - PCRE: perl=TRUE option



# Outside of the data box

## Change the problems

---

- Regression to Classification (vice versa)
  - Numeric variables to Characteristic variables (vice versa)
  - Word data(categorical) to Document(Text)
  - And many more...
-



# Purchase history data



## But...Data?

---

- Fried Chicken (치킨)
  - 후라이드치킨
  - 크리스피치킨
  - 후라이드 치킨
  - 마일드치킨(Mild)
  - 오리지널치킨
  - 닭강정★
  - 써프라이드 - 매운맛
  - 소이갈릭스
  - 1. 치즐링
  - A. 순살후라이드
  - 마라치킨(麻辣)
  - ....

# Restaurants and Users

## Restaurants

- Each restaurant sells various dishes



## Users

- Every person eats various food



- Food name cleansing with regular expressions and rules

# Restaurant Documents

- [illegible]

# Text Processing

## Tidyttext

---

- Text mining for word processing and sentiment analysis using 'dplyr', 'ggplot2', and other tidy tools.

## Example

---

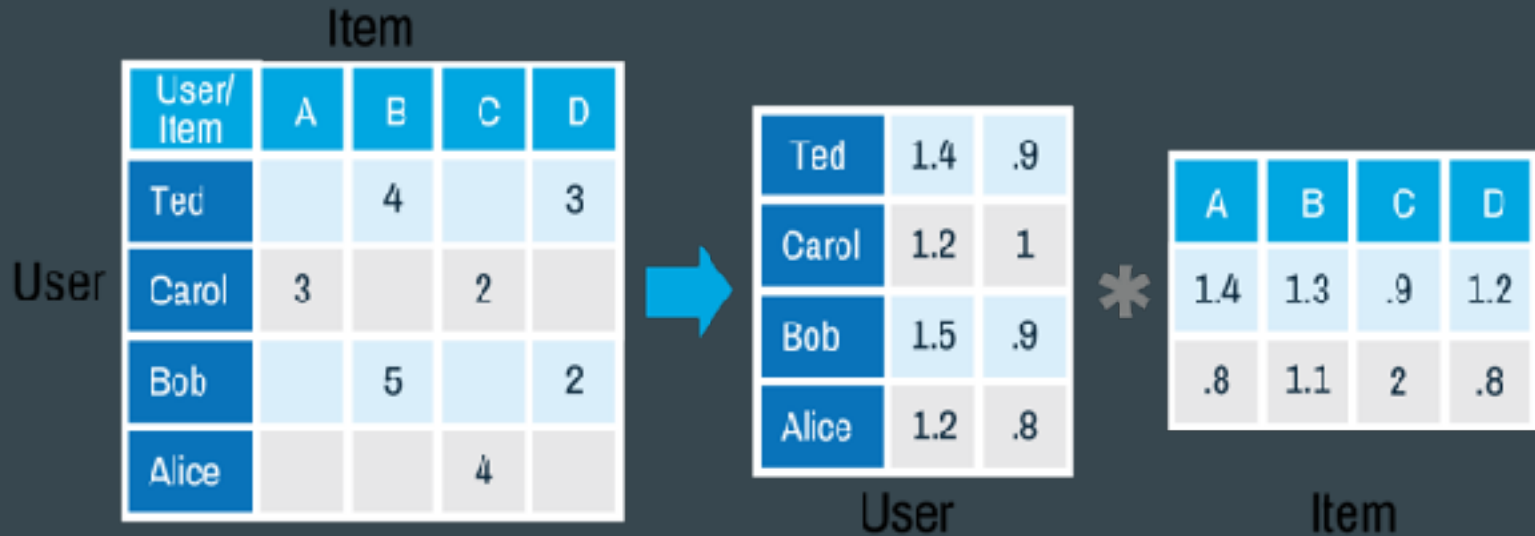
```
library(tidyverse)

shopTF <- shopDoc %>% unnest_tokens(word, shop_food_doc)
               %>% count(theme_seq, shop_no, word)
               %>% ungroup()
shopTF1 <- cbind(paste(shopTF$theme_seq, shopTF$shop_no),
                 shopTF[,c(3,4)])
colnames(shopTF1) <- c("theme_shop", "word", "n")
```

# Recommendation with TF-IDF

## ALS algorithms

- Users and restaurants(Item) matching



# Summary

## Machine Learning Data

There are far fewer clean data for machine learning, but we can make with data wrangling

## Data Wrangling

Tidyverse + Regular Expression + Out of the box thinking



The Era of **Data Wrangling**  
is coming....



# Q & A

Thank you for listening

Contact: [cojette@gmail.com](mailto:cojette@gmail.com)

@mazycat