



# The life-changing magic of tidying up

---

*Min Fang*



“

“80% of data analysis is spent on  
the process of cleaning and  
preparing the data.”

— *Conventional wisdom*

# Tidy data

“

Each variable forms a column.

Each observation forms a row.

Each type of observational unit forms a table.

— Hadley Wickham

# R packages I use regularly to handle tidy data seamlessly

---

- tidyverse (<https://www.tidyverse.org/>)
  - readr
  - tidyverse
  - dplyr
  - ggplot2
- stringr

**Doing things easily with  
tidy data**

# Q1: What is the distribution of aspect labels from the model?

```
In [6]: 1 head(sentiment_eval)
```

# Q1: What is the distribution of aspect labels from the model?

---

```
In [21]: 1 predictions <- sentiment_eval %>% select(-starts_with("QA_"))
2 head(predictions)
```

review_id	review	bathroom	breakfast	internet	staff	swimming_pool
1	Was ment to stay for 2 nights lasted one! The bath [...]	3	NA	NA	NA	NA
2	On arrival walking to reception doors met by cig e [...]	3	NA	NA	1	NA
3	Had read some reviews before my visit on 17th Oct' [...]	NA	NA	NA	1	3
4	Went with open mind after reading recent reviews b [...]	3	3	NA	3	NA
5	Expensive rooms did not match expectations. Breakf [...]	NA	NA	NA	NA	NA
6	Overall, great hotel with big rooms and toilets, g [...]	3	NA	NA	NA	NA

# Q1: What is the distribution of aspect labels from the model?

---

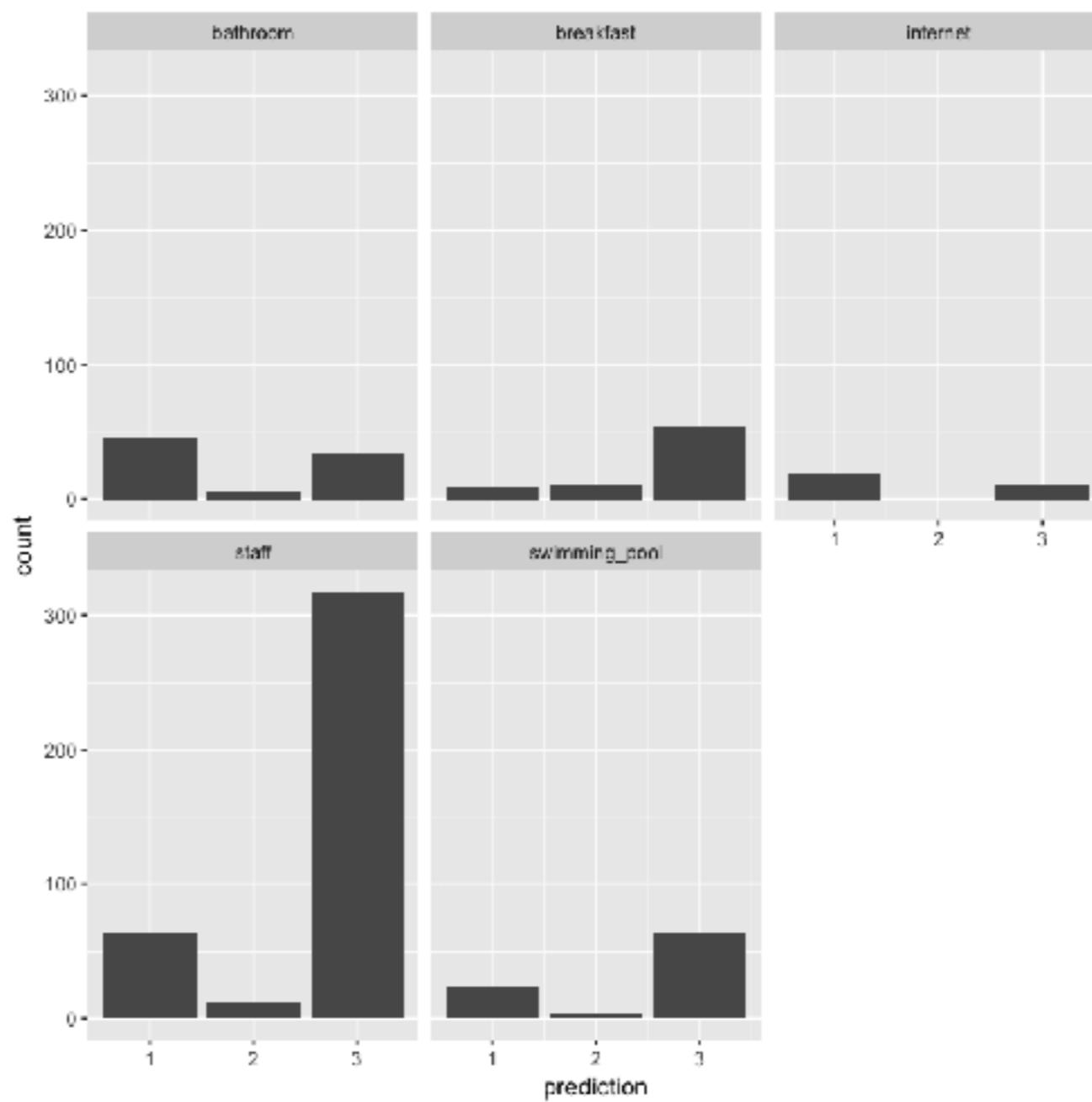
```
In [25]: 1 sentiment_eval %>% select(-starts_with("QA_")) %>%
2   gather("aspect", "prediction", 3:7) %>% arrange(review_id)
```

review_id	review	aspect	prediction
1	Was ment to stay for 2 nights lasted one! The bath [...]	bathroom	3
1	Was ment to stay for 2 nights lasted one! The bath [...]	breakfast	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	internet	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	staff	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	swimming_pool	NA
2	On arrival walking to reception doors met by cig e [...]	bathroom	3
2	On arrival walking to reception doors met by cig e [...]	breakfast	NA
2	On arrival walking to reception doors met by cig e [...]	internet	NA
2	On arrival walking to reception doors met by cig e [...]	staff	1
2	On arrival walking to reception doors met by cig e [...]	swimming_pool	NA
3	Had read some reviews before my visit on 17th Oct' [...]	bathroom	NA
3	Had read some reviews before my visit on 17th Oct' [...]	breakfast	NA
3	Had read some reviews before my visit on 17th Oct' [...]	internet	NA
3	Had read some reviews before my visit on 17th Oct' [...]	staff	1
3	Had read some reviews before my visit on 17th Oct' [...]	swimming_pool	3
4	Went with open mind after reading recent reviews b [...]	bathroom	3
4	Went with open mind after reading recent reviews b [...]	breakfast	3
4	Went with open mind after reading recent reviews b [...]	internet	NA
4	Went with open mind after reading recent reviews b [...]	staff	3

# Q1: What is the distribution of aspect labels?

---

```
1 predictions.tidy <- sentiment_eval %>% select(-starts_with("QA_")) %>%
2   gather("aspect", "prediction", 3:7) %>% arrange(review_id)
3 ggplot(predictions.tidy) + facet_wrap(~aspect) +
4   geom_bar(aes(x=prediction, y=..count..))
```



## Q2: What is the number of aspects labels\* per review?

---

```
In [25]: 1 sentiment_eval %>% select(-starts_with("QA_")) %>%
2   gather("aspect", "prediction", 3:7) %>% arrange(review_id)
```

review_id	review	aspect	prediction
1	Was ment to stay for 2 nights lasted one! The bath [...]	bathroom	3
1	Was ment to stay for 2 nights lasted one! The bath [...]	breakfast	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	internet	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	staff	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	swimming_pool	NA
2	On arrival walking to reception doors met by cig e [...]	bathroom	3
2	On arrival walking to reception doors met by cig e [...]	breakfast	NA
2	On arrival walking to reception doors met by cig e [...]	internet	NA
2	On arrival walking to reception doors met by cig e [...]	staff	1
2	On arrival walking to reception doors met by cig e [...]	swimming_pool	NA
3	Had read some reviews before my visit on 17th Oct' [...]	bathroom	NA
3	Had read some reviews before my visit on 17th Oct' [...]	breakfast	NA
3	Had read some reviews before my visit on 17th Oct' [...]	internet	NA
3	Had read some reviews before my visit on 17th Oct' [...]	staff	1
3	Had read some reviews before my visit on 17th Oct' [...]	swimming_pool	3
4	Went with open mind after reading recent reviews b [...]	bathroom	3
4	Went with open mind after reading recent reviews b [...]	breakfast	3
4	Went with open mind after reading recent reviews b [...]	internet	NA
4	Went with open mind after reading recent reviews b [...]	staff	3

\* or any other interesting aggregate function

## Q2: What is the number of aspects labels\* per review?

---

```
1 predictions.tidy %>% group_by(review) %>%
2 summarise(
3   num_predictions = sum(!is.na(prediction))
4 )
```

	review	num_predictions
1	1st my wife asked for help with lugagges and told [...]	2
2	A very good hotel with a wonderful pool,, excellen [...]	2
3	After our problem with check-in our stay was quite [...]	1
4	After three weeks in Vietnam, spending a few days [...]	2
5	Air conditioners in the rooms are very loud other [...]	1
6	All fine, its a huge hotel really nice location, v [...]	1
7	All I can say is "WOW"....the staff is amazing!! [...]	1
8	All of the staff at this hotel were amazingly frie [...]	1
9	All staff friendly & helpful. Room was comfortable [...]	1
10	Almost everything about this place is 5 star. Our [...]	1
11	Amazing grounds. Amenities and staff were amazing. [...]	1
12	An error was made on the 2 rooms booked in but the [...]	0
13	An old property, but the pool is adequate, the bed [...]	2
14	Animal friendly, staff was friendly and helpful. W [...]	1

## Q2: What is the number of aspects labels\* per review?

---

```
1 predictions.tidy %>% group_by(review) %>%
2 summarise(
3   num_predictions = sum(!is.na(prediction)),
4   max_prediction = max(prediction, na.rm = T),
5   min_prediction = min(prediction, na.rm = T),
6   avg_prediction = mean(prediction, na.rm = T),
7   max_min_diff = max_prediction - min_prediction,
8   sd_prediction = sd(prediction, na.rm = T)
9 )
```

review	num_predictions	max_prediction	min_prediction	avg_prediction	max_min_diff	sd_prediction
1st my wife asked for help with lugagges and told [...]	2	3	3	3.0	0	0.0000000
A very good hotel with a wonderful pool,, excellen [...]	2	3	3	3.0	0	0.0000000
After our problem with check-in our stay was quite [...]	1	3	3	3.0	0	NA
After three weeks in Vietnam, spending a few days [...]	2	3	2	2.5	1	0.7071068
Air conditioners in the rooms are very loud other [...]	1	2	2	2.0	0	NA
All fine, its a huge hotel really nice location, v [...]	1	3	3	3.0	0	NA
All I can say is "WOW"....the staff is amazing!! [...]	1	3	3	3.0	0	NA
All of the staff at this hotel were amazingly frie [...]	1	3	3	3.0	0	NA
All staff friendly & helpful. Room was comfortable [...]	1	3	3	3.0	0	NA
Almost everything about this place is 5 star. Our [...]	1	1	1	1.0	0	NA
Amazing grounds. Amenities and staff were amazing. [...]	1	3	3	3.0	0	NA
An error was made on the 2 rooms booked in but the [...]	0	-Inf	Inf	NaN	-Inf	NA
An old property, but the pool is adequate, the bed [...]	2	3	3	3.0	0	0.0000000

# Q3: How much do the two opinions agree per aspect?

```
In [6]: 1 head(sentiment_eval)
```

# Q3: How much do the two opinions agree per aspect?

---

review_id	review	aspect	prediction	qa
1	Was ment to stay for 2 nights lasted one! The bath [...]	bathroom	3	2
1	Was ment to stay for 2 nights lasted one! The bath [...]	breakfast	NA	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	internet	NA	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	staff	NA	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	swimming_pool	NA	NA
2	On arrival walking to reception doors met by cig e [...]	bathroom	3	2
2	On arrival walking to reception doors met by cig e [...]	breakfast	NA	NA
2	On arrival walking to reception doors met by cig e [...]	internet	NA	NA
2	On arrival walking to reception doors met by cig e [...]	staff	1	2
2	On arrival walking to reception doors met by cig e [...]	swimming_pool	NA	NA
3	Had read some reviews before my visit on 17th Oct' [...]	bathroom	NA	NA
3	Had read some reviews before my visit on 17th Oct' [...]	breakfast	NA	NA
3	Had read some reviews before my visit on 17th Oct' [...]	internet	NA	NA
3	Had read some reviews before my visit on 17th Oct' [...]	staff	1	3
3	Had read some reviews before my visit on 17th Oct' [...]	swimming_pool	3	3

Goal dataframe

# Q3: How much do the two opinions agree per aspect?

review_id	review	aspect	prediction
1	Was ment to stay for 2 nights lasted one! The bath [...]	bathroom	3
1	Was ment to stay for 2 nights lasted one! The bath [...]	breakfast	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	internet	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	staff	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	swimming_pool	NA
2	On arrival walking to reception doors met by cig e [...]	bathroom	3

```
1 qa.tidy <- sentiment_eval %>% select(review_id, review, starts_with("QA_")) %>%
2 gather("aspect", "qa", 3:7) %>% arrange(review_id)
3 head(qa.tidy)
```

review_id	review	aspect	qa
1	Was ment to stay for 2 nights lasted one! The bath [...]	QA_bathroom	2
1	Was ment to stay for 2 nights lasted one! The bath [...]	QA_breakfast	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	QA_internet	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	QA_staff	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	QA_swimming_pool	NA
2	On arrival walking to reception doors met by cig e [...]	QA_bathroom	2

# Q3: How much do the two opinions agree per aspect?

---

review_id	review	aspect	prediction
1	Was ment to stay for 2 nights lasted one! The bath [...]	bathroom	3
1	Was ment to stay for 2 nights lasted one! The bath [...]	breakfast	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	internet	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	staff	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	swimming_pool	NA
2	On arrival walking to reception doors met by cig e [...]	bathroom	3

```
1 qa.tidy <- qa.tidy %>% mutate(aspect = str_replace(aspect, "QA_", ""))
2 head(qa.tidy)
```

review_id	review	aspect	qa
1	Was ment to stay for 2 nights lasted one! The bath [...]	bathroom	2
1	Was ment to stay for 2 nights lasted one! The bath [...]	breakfast	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	internet	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	staff	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	swimming_pool	NA
2	On arrival walking to reception doors met by cig e [...]	bathroom	2

# Q3: How much do the two opinions agree per aspect?

---

```
1 evaluation <- inner_join(predictions.tidy, qa.tidy, by=c('review_id', 'review', 'aspect'))  
2 evaluation
```

review_id	review	aspect	prediction	qa
1	Was ment to stay for 2 nights lasted one! The bath [...]	bathroom	3	2
1	Was ment to stay for 2 nights lasted one! The bath [...]	breakfast	NA	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	internet	NA	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	staff	NA	NA
1	Was ment to stay for 2 nights lasted one! The bath [...]	swimming_pool	NA	NA
2	On arrival walking to reception doors met by cig e [...]	bathroom	3	2
2	On arrival walking to reception doors met by cig e [...]	breakfast	NA	NA
2	On arrival walking to reception doors met by cig e [...]	internet	NA	NA
2	On arrival walking to reception doors met by cig e [...]	staff	1	2
2	On arrival walking to reception doors met by cig e [...]	swimming_pool	NA	NA
3	Had read some reviews before my visit on 17th Oct' [...]	bathroom	NA	NA
3	Had read some reviews before my visit on 17th Oct' [...]	breakfast	NA	NA
3	Had read some reviews before my visit on 17th Oct' [...]	internet	NA	NA
3	Had read some reviews before my visit on 17th Oct' [...]	staff	1	3
3	Had read some reviews before my visit on 17th Oct' [...]	swimming_pool	3	3

# Q3: How much do the two opinions agree per aspect?

---

```
1 evaluation %>% mutate(agree = prediction == qa %>%
2 group_by(aspect, agree) %>% summarise(
3   count = n()
4 )
```

aspect	agree	count
bathroom	FALSE	20
bathroom	TRUE	15
bathroom	NA	535
breakfast	FALSE	6
breakfast	TRUE	37
breakfast	NA	527
internet	FALSE	4
internet	TRUE	11
internet	NA	555
staff	FALSE	36
staff	TRUE	172
staff	NA	362
swimming_pool	FALSE	9
swimming_pool	TRUE	41
swimming_pool	NA	520

```
1 evaluation %>% filter(!is.na(prediction) | !is.na(qa)) %>%
2 mutate(agree = if_else(
3   is.na(prediction == qa), F, prediction == qa
4   )
5 ) %>%
6 group_by(aspect, agree) %>% summarise(
7   count = n()
8 )
```

aspect	agree	count
bathroom	FALSE	90
bathroom	TRUE	15
breakfast	FALSE	49
breakfast	TRUE	37
internet	FALSE	27
internet	TRUE	11
staff	FALSE	249
staff	TRUE	172
swimming_pool	FALSE	65
swimming_pool	TRUE	41

# Q3: How much do the two opinions agree per aspect?

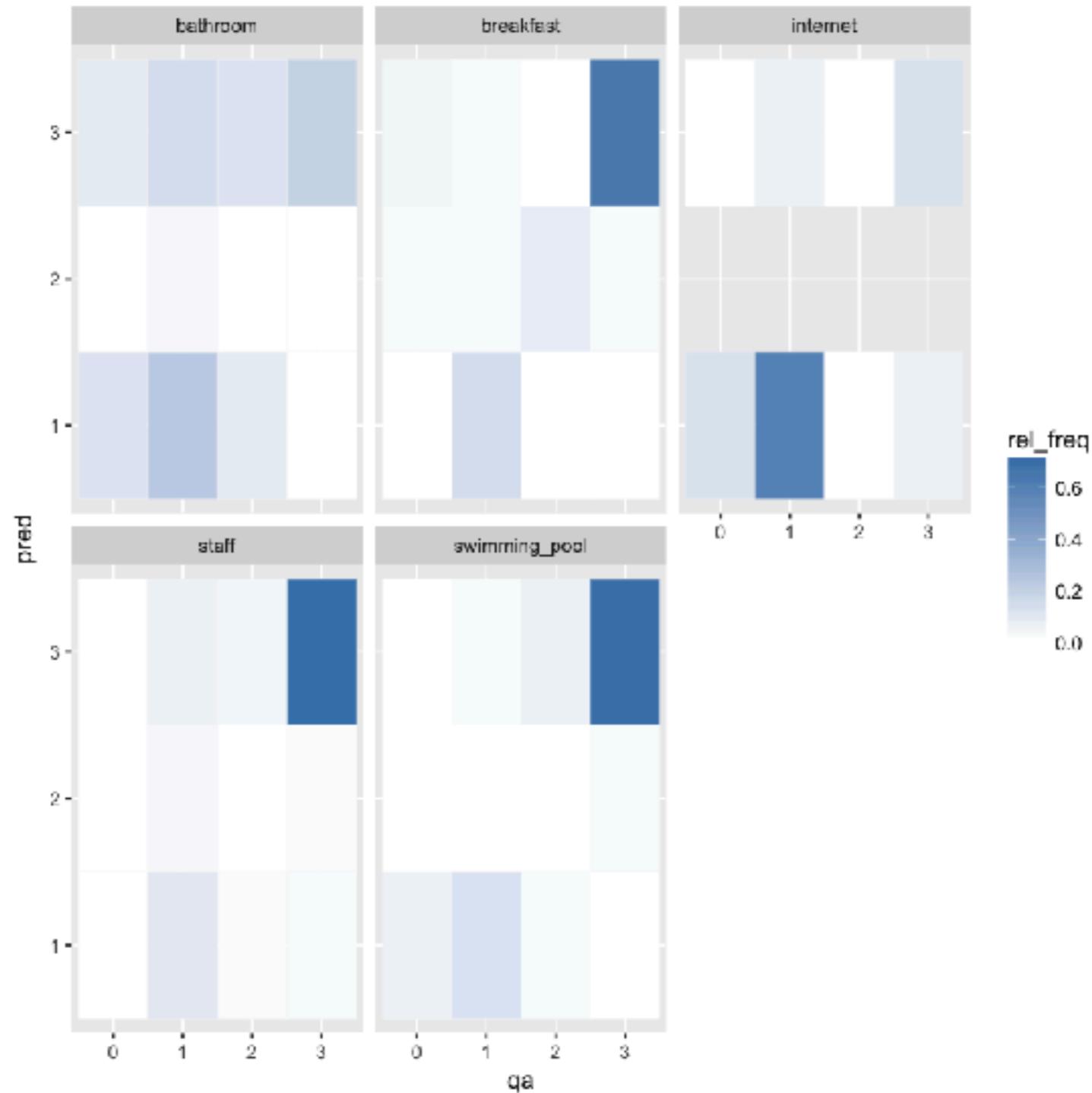
---

```
1 evaluation %>% mutate(agree = case_when(  
2     is.na(prediction) & is.na(qa) ~ T,  
3     is.na(prediction == qa) ~ F,  
4     !is.na(prediction == qa) ~ prediction == qa  
5 )) %>%  
6 group_by(aspect, agree) %>% summarise(  
7     count = n()  
8 ) %>% mutate(rel_freq = count/sum(count))
```

aspect	agree	count	rel_freq
bathroom	FALSE	90	0.15789474
bathroom	TRUE	480	0.84210526
breakfast	FALSE	49	0.08596491
breakfast	TRUE	521	0.91403509
internet	FALSE	27	0.04736842
internet	TRUE	543	0.95263158
staff	FALSE	249	0.43684211
staff	TRUE	321	0.56315789
swimming_pool	FALSE	65	0.11403509
swimming_pool	TRUE	505	0.88596491

# Q3: How much do the two opinions agree per aspect?

---



# Q3: How much do the two opinions agree per aspect?

---

```
1 crosstab <- evaluation %>% group_by(aspect) %>% do(  
2   data_frame(crosstab=list(table(qa=.qa, pred=.prediction)))  
3 )  
4 crosstab
```

aspect	crosstab
bathroom	4, 8, 3, 0, 0, 1, 0, 0, 3, 5, 4, 7
breakfast	0, 6, 0, 0, 1, 1, 3, 1, 2, 1, 0, 28
internet	2, 9, 0, 1, 0, 1, 0, 2
staff	0, 21, 3, 5, 0, 6, 0, 3, 0, 11, 8, 151
swimming_pool	3, 6, 1, 0, 0, 0, 0, 1, 0, 1, 3, 35

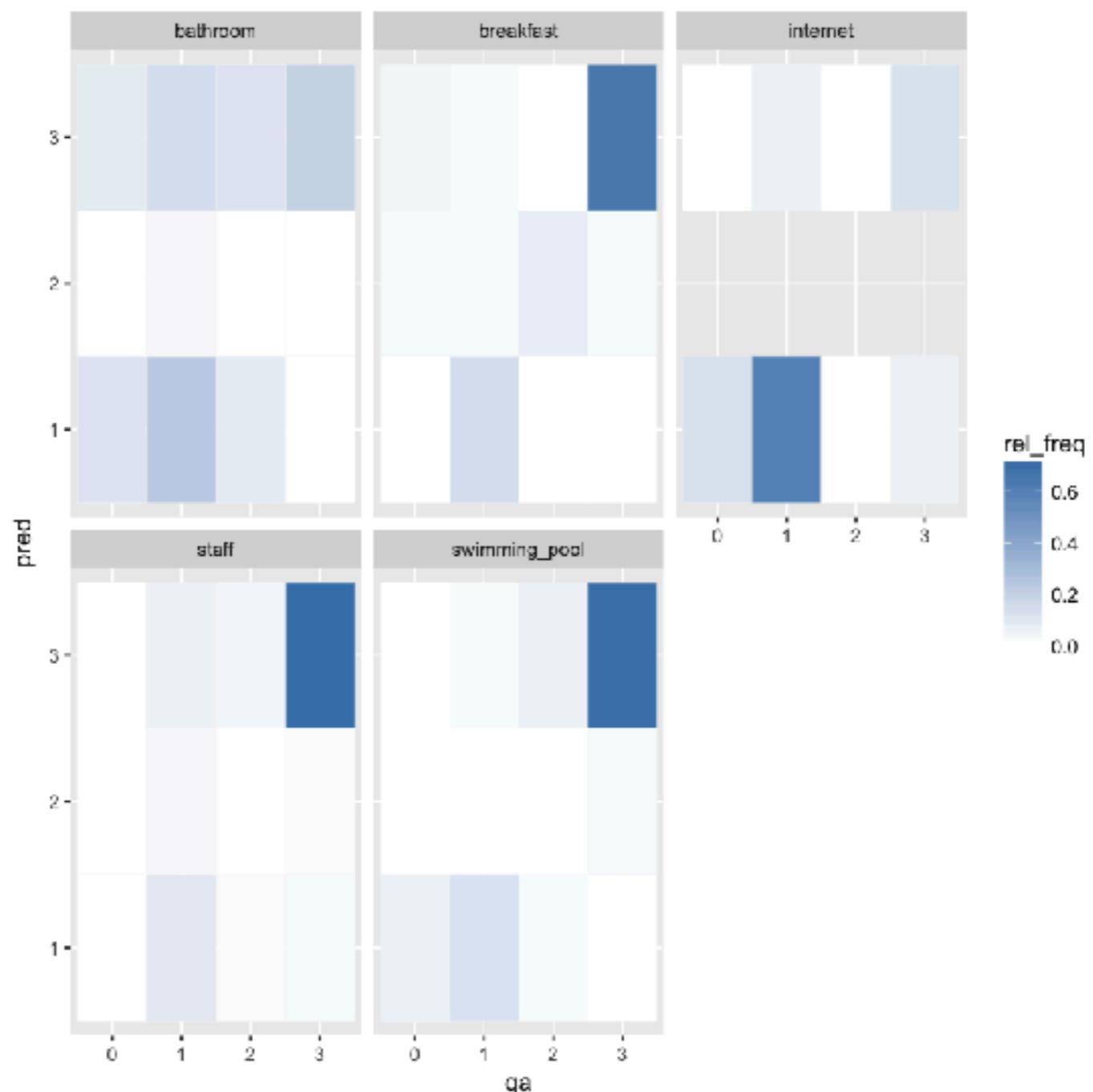
```
1 crosstab.tidy <- crosstab %>% group_by(aspect) %>%  
2 do(as.data.frame(.scrosstab)) %>%  
3 group_by(aspect) %>% mutate(rel_freq = Freq/sum(Freq))  
4 crosstab.tidy
```

aspect	qa	pred	Freq	rel_freq
bathroom	0	1	4	0.11428571
bathroom	1	1	8	0.22857143
bathroom	2	1	3	0.08571429
bathroom	3	1	0	0.00000000
bathroom	0	2	0	0.00000000
bathroom	1	2	1	0.02857143
bathroom	2	2	0	0.00000000
bathroom	3	2	0	0.00000000
bathroom	0	3	3	0.08571429
bathroom	1	3	5	0.14285714
bathroom	2	3	4	0.11428571
bathroom	3	3	7	0.20000000
breakfast	0	1	0	0.00000000

# Q3: How much do the two opinions agree per aspect?

---

```
1 ggplot(crosstab.tidy, aes(qa, pred)) +
2   facet_wrap(~aspect) +
3   geom_tile(aes(fill = rel_freq), colour = "white") +
4   scale_fill_gradient(low = "white", high = "steelblue")
```



**Textual data can be  
tidy too!**

**spa**  
friendly service staff positive  
clean quality located  
amazing bath bar restaurant  
close people free sauna morning massage  
times enjoyed views bit stayed city lot wonderful  
price walking quiet top beautiful front hotels bed excellent  
spacious return trip beach wifi dinner star nights booked  
check beds resort town guests buffet  
feel food nice stay loved main kids  
minutes hot evening day visit facilities staying water  
fantastic gym night walk coffee floor  
perfect relaxing pools easy view staying  
experience 5 days highly pool  
comfortable helpful family  
recommend bathroom  
**breakfast**  
restaurants

spa excellent lovely surroundings  
water villas relations manager  
**nice spa** amazing helpful services negative  
friendly spa incredibly accommodating  
swedish massage spa fantastic  
relaxing massage spa perfect fantastic spa  
class spa relaxing spa hotel spa  
signature massage spa negative spa amazing  
brilliant spa spa negative onsite spa  
spa offers relaxing getaway spa amazing  
location amazing comfy clean free spa balinese massage  
treatments offered location spa food spa amazing  
inclusive hotels enjoyed spa massive bed  
complete relaxation spa massages spa incredible spa  
superb spa liva spa body massage positive pool  
positive spa superb spa liva wellness spa spa wonderful  
beautiful spa spa treatment spa services  
site spa spa steam spa treatments  
free beach traditional service minded coconut trees  
lovely spa tradafabulous friendly super spa  
hammam treatment body massagess quiet beautiful  
enchantment resort perfect spa reception helped  
awesome pool wonderful spa wife loved day spa  
body lotion spa experience service lovely michelin star  
pacific hotel boutique resort breakfast everyday warm staff  
body scrub excellent spa spa lovely fabulous spa  
excellent massage spa facilities stone massage  
kids enjoyed spa nice fabulous massage  
wonderful massage excellent gym  
aromatherapy massage  
fantastic massage  
  
**negative** helpful negative  
extremely disappointed  
black hair hotel massage  
shower screen negative extra  
**positive** negative  
**negative spa**  
negative pool day negative super cheap art deco  
art deco  
luke warm booked 4 day 2 terrible experience positive  
lovely city hotel advertises price negative  
spa appointment negative swimming  
positive comfortable wedding venue normal price  
oil massage positive location facilities beautiful  
friendly negative spa prices ice bucket xuan huong  
brilliant negative positive staff wash basin  
desk persongym steam poor experience bed minute wait fire alarm  
secret escapes 00 pm kerala massage hotel 3 terrible service  
plan aheadbar negative rock hard 5.7 worst experience night due staff told costs 10  
worst experience mins massage rated hotel bit dirty  
bed comfy positive beds spa staff change towels low quality  
singles pushed previous guest star quality 7 pm  
freezing cold north indian bit disappointing superior double  
45 minute heritage property 5.00 answer questions  
people talking fairly close bit awkward 2 hotels held shower  
spa section biggest disappointment week ago  
positive spacious positive food negative 1  
comfortable location level 2  
stay negative spacious negative ladies changing  
comfortable negative cold breakfast  
negative sauna  
location negative  
spa closed

# References & Resources

---

- <http://vita.had.co.nz/papers/tidy-data.pdf>
- <http://tidyr.tidyverse.org/articles/tidy-data.html>
- Best cheatsheet ever! <https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf>
- tidytext book: <https://www.tidytextmining.com/>
- widyr: <https://github.com/dgrtwo/widyr>
- dplyr for Python: <https://github.com/kieferk/dfply>
- Pandas ‘equivalents’ dplyr: [https://pandas.pydata.org/pandas-docs/stable/comparison\\_with\\_r.html](https://pandas.pydata.org/pandas-docs/stable/comparison_with_r.html)



Min Fang

Data Scientist in Hotel Profiling  
trivago Amsterdam

Email: [min.fang@trivago.com](mailto:min.fang@trivago.com)

