

### Trifacta Cleaning

1. cleaned column with all NA or empty or just several useless values: thumbnail\_url, medium\_url, xl\_picture\_url, neighbourhood\_group\_cleansed, experiences\_offered, host\_acceptance\_rate, license, jurisdiction\_names
2. cleaned column with meaningless same value throughout the sheet: scrape\_id, neighbourhood, city, market, smart\_location, country\_code, street, host\_location, host\_neighbourhood, host\_listings\_count, calendar\_last\_scraped, has\_availability, requires\_license, is\_business\_travel\_ready
3. delete all text description and pictures( we might add them back in next stages analysis if we decide to go fancy and applying some NLP, for no we only look at those numeric variables: name, summary, space, description, neighborhood\_overview, notes, transit, access, interaction, house\_rules, host\_about
4. Delete all urls which we might also add back later but for now just useless: picture\_url, host\_thumbnail\_url, host\_picture\_url
5. Delete listing\_url(= [https://www.airbnb.com/rooms/listing\\_id](https://www.airbnb.com/rooms/listing_id)), host\_url(=https://www.airbnb.com/users/show/host\_id)totally meaningless.
6. Although we only leave numbers and deleted both, just a reminder, summary and description are very similar column.
7. Change the data type of "reply rate" to numbers and replace N/A to average value.
8. Make all empty or 0 value in host\_total\_listings\_count / square\_feet /amenities/square\_feet as NULL
9. Delete "" and space in the amenities text and host\_verifications
10. Clean the empty and "" value in zip code to NULL, for some zip code longer than 5 digit, only leave the first five digit, all Boston area zip code should start with 0, make those wired zip code NULL as well.
11. I don think we care about the maximum\_night we should only care about the minimum night at this case( if I am wrong we can add it back anyway, but for now let's just delete other column.
12. Calculate the average price by host and add as a new column by left join.
13. Still looking for a way to turn amenities and host\_verifications into factors but I think that's pretty much good for now.