

# Multi-Doc Summarization

## COVID-19

Sharon Wang  
Week 5

# CatE

Performance of CatE depends on:

1. **corpus context statistics**; if a category name only appears a few times in the corpus, there may not be good contexts for the category name
2. **input category names**; if a category name doesn't have other category names to be distinguished from, the results may be too general
  - a. for category names China, Italy, United\_States, and New\_Zealand, the output for Italy is: ->
  - b. **solution: add other countries (e.g. Spain, France, Germany, ...) to the list of category names**
  - c. alternative solution: discard the other country names from the output

<b>Italy</b>
Spain
Lombardy
Europe
Italys
France
Germany
Milan
Venice
Veneto
Italy

# Corpus Context Statistics

1. Count the number of **declarative sentences** containing
  - a. the category name
  - b. a **subject**, **predicate**, and **object** (using SRL)
2. Only run CatE on the category name if the number of declarative sentences > ~50.

# Input Category Names

- To get more specific results for category names, use SpaCy and WordNet to find other words that are in the same classification as the category name.
- For example, for the category name “Italy,” use SpaCy to classify it as a country.
- Then use WordNet to generate a certain number of other words that are both countries and appear in the corpus.
- **Alternatively, can go through the corpus and add any words classified as a country to the category name list.**
- Then add these words as other category names and run CatE for a better result.

# Results

- Applying these new rules, the new category names will be close to: **China, Italy, United\_States, Spain, France, Germany, Iran, Indonesia, Switzerland, etc.**
- Notes:
- New\_Zealand was removed because there were only 9 declarative sentences containing “New\_Zealand”, and a subject, predicate, and object in the corpus.
- Spain, France, and Germany, etc were added because they are other countries and they appear in the corpus in quality sentences more than 50 times.

# CatE Category Names (Countries)

```
1 Category (China):
2 Wuhan mainland Chinese Hubei Chinas Li human-to-human China Shanghai Beijing
3 Category (Italy):
4 Tuscany happen in ITALYThe Christie Korea northern Italys Lombardy Stettler towns
5 Category (United_States):
6 recent_days 41000 53000 confirmed_coronavirus_cases surpassed topped surged exceeds total_infections neared
7 Category (Spain):
8 hard-hit devastation hardest-hit morgues Tamlin overwhelmed overwhelming hallways Maisonneuve-Rosemont hospitals
9 Category (France):
10 Bishops French Paris QESBA parents lessons English Kingdom courses exams
11 Category (Germany):
12 Daniele Rome Italian Alberto Spain Milan Madrid Spains wheeled Ho
13 Category (Iran):
14 Europe Egypt Italy Milan Korea Spain Asia Italys Lombardy Italian
15 Category (Indonesia):
16 Willy KurniawanReuters worlds Smithfield_Foods biggest giant halted second-largest tankers company said
17 Category (Switzerland):
18 Anutin Bordeaux Charnvirakul returned from reporters Speaking Dominic Networks helpful spoke
19 Category (Ukraine):
20 Op Croizer recorders Unifier Defence authorization recorder Navys Capt helmet
21 Category (Greece):
22 military NATO troops peacekeeping mission soldiers Forces helicopter reservists fighter
23 Category (Turkey):
24 returned 27-year-old who travelled woman who returning travelled 14-day internationally flew recently
25 Category (Algeria):
26 confirmed in cruise to travel BC – Two Three March Princess cruise 50s who three
27 Category (Lithuania):
28 Oriental Huntington Palm Orpheum Finals Boutique Mapelli Mozzi Alpine Ski
29 Category (Slovakia):
30 shopping_malls showrooms shops close restaurants retail movie theatres venues spas
31 Category (Denmark):
32 gradual_reopening public_transport phased-in restrictions to approach workplaces loosening gradual – such province – if
33 Category (Australia):
34 Jacinda Bahrain headline Zealand tighten bother Medicare List Ardenr favouring
35 Category (Egypt):
36 Iran COVID-19 in whereabouts coronavirus in asked self-isolate Iran case Ayatollah Shahab
37 Category (Sweden):
38 Asian_countries wallow squandered watchdogs critics stoked likelier stereotypes mid-January inclined
39 Category (South_Africa):
40 strategy expressed_concern highlighting diversified economy creating wrought prop fast-growing broadest
41 Category (Slovenia):
42 Cyprus border controls barred traffic restrict crossing Canada-US crossings entry
43 Category (Austria):
44 bakeries supermarkets stores pharmacies grocery soups foods hardware sauce meats
45 Category (Argentina):
46 Tucumán Miróli Schaffner Américo immunologist pathology Network associate_professor Alfredo Susy
47 Category (Ghana):
48 727 three new 405 its total also Tuesday On Wednesday Thursday afternoon
```

- The countries we are testing is **China**, **Italy**, and **United\_States**. The other category names are included so they do not show up in the first 3 countries' generated phrases. (There are even more countries, but they couldn't fit onto the slide).
- This improved the output, but it's still not perfect.
- For example, "Korea" was generated for Italy, even though it was included as a category name.
- However, the error is much lower now so as a solution, results that are in the same classification as the category names could be deleted? I.e., delete "Korea" from Italy's phrases since Korea is a country.