Data for "Gender Bias in Rumors among Professionals: An Identity-based Interpretation"

### Alice Wu

This file lists all data sets provided for replication in the "data" folder. Please also refer to "README code.pdf" for programs that generate and clean the datasets.

### > Final data sets:

- "full\_sample2019\_stata\_final.dta": all threads in the gender sample that include at least one Female/Male post. Variables include measures of gender and topics at the posts level. Please see "codebook\_full\_sample2019\_dta.log" for variable descriptions. All results in this paper are run on this dataset, through "static\_analysis\_2019.do" and "dynamic\_analysis\_2019.do"
- "full\_sample2019\_stata.csv": csv format for the dta file above
- "full\_sample2019.csv": all threads in gender sample. The column "raw\_cat" contains the content of original posts. Please see "codebook\_datasets.xlsx" for variable descriptions. The step from this dataset to the final dataset "full\_sample2019\_stat.csv" is completed by "gen\_full\_sample.R".

Now I start from the raw data scraped from the Economics Job Market Rumors forum and explain how I generate the final datasets above.

### Raw data sets:

- (a) Directory: "..data/raw"
  - "scraped\_posts.txt": raw text for threads scraped from the EJMR forum as of October 2017
  - "EJRO\_raw\_text.csv": scraped data above converted into a .csv file
  - "EJRO\_raw\_text\_cleaned.csv": raw data with unique IDs for threads, and posts under each thread, cleaned from "EJRO\_raw\_text.csv" through the program "clean\_raw\_text.R"
  - "main\_stats\_cleaned.csv": thread-level statistics as shown on the main listings on the EJMR websites, including information such as thread title, total number of posts and views, etc.
  - **"raw\_time\_stamp.csv"**: time stamp (e.g., "Posted on: 6 years ago") for the first post in each thread
  - "vocab\_Nov2017.pkl": a dictionary for the most frequent 10,000 words, generated through the program "word\_count.py"

- "counts\_Nov2017.npz": 2.2 million by 10,000 matrix of word counts saved in compressed format, also generated through "word\_count.py".
- "all\_cat\_2019.csv": counts of words under each category in each post. The category assignment for each of the most frequent 10,000 words can be found at "../vocab/cleaned\_vocab.csv".
- "EJR\_ALL\_categories\_2019.csv": same as above but include IDs for threads and posts, see "prep\_for\_analysis.py".
- "all\_classifiers\_2019.csv": counts of female and male classifiers in each post.
- **"EJR\_ALL\_gender\_classifiers\_2019.csv"**: same as above but include IDs for threads and posts, see "prep\_for\_analysis.py".
- "full\_sample\_job\_rank.csv": counts of keywords indicating each job rank in each post in the full gender sample (threads that include at least one Female or Male post), generated by "job-rank.R".

# (b) Directory: "..data/vocab"

- "cleaned\_vocab.csv": the most frequent 10,000 words as in "vocab\_Nov2017.pkl", indexed from 1 to 10,000, and category classification
- **"gender\_classifiers.csv"**: 53 female and 204 male classifiers, e.g., gender pronouns, woman/man.
- "category" folder: each .txt file here lists the indices under each category

## (c) <u>Directory: "../data/gender\_sample\_by\_classifiers"</u>

 "EJR\_gender\_dataset\_Jan2018.csv": threads that include at least one Female or Male post identified through gender classifiers alone (not names of economists as of January 2018). This sample includes 1.7 million posts in total. The column "female0\_pred" indicates if it is a Female versus Male post.

## (d) Directory: "..data/NBER"

- "author-history-ID-by-full.csv": thread-post IDs for posts that mention an NBER author's full name, generated by "NBER-post-history-ID (by full).R"
- "author-history-ID-by-part.csv": IDs for posts that mention first name/last name of an author, within threads that include at least one post mentioning his or her full name; generated by "NBER-post-history-ID (by part).R"

- "author-history-merged.csv": merge the datasets above with 2.2 million posts in "EJRO\_raw\_text\_cleaned.csv", and preserve all posts under threads with at least one occurrence of an author's full name.
- "nber-author-info.csv": information about authors of NBER working papers, scraped from the NBER websites.
- "complete-nber-author-info.csv": information about NBER authors as above with manually assigned job rank (when one is not affiliated with NBER).

## (e) Directory: "..data/JMC"

- "jmc\_data\_gender\_nonmissing.csv": a directory for 4,750 economics Ph.D. who graduated between 2011 and 2018. The dataset includes their names, gender, school and year of graduation.
- "JMC-history-ID.csv": thread-post IDs for posts that include a person's full name, or first name/last name/initials under threads where full name has occurred at least once. The dataset is generated by "JMC-post-history-ID.R".
- "JMC-history-ID-merged.csv": merge the two datasets above.

(c) (d) (e) provide three sources of posts that discuss female or male subjects. I merge them together and generate the full gender sample in "full\_sample2019.csv" that preserve all posts under threads that include at least one Female or Male post.