# Toward a semi-automated item nonresponse detector model for open-response data

Kristen Cibelli Hibben, PhD; Zachary Smith, MA; Travis Hoppe, PhD; Valerie Ryan, PhD; Ben Rogers, MS; Paul Scanlon, PhD; Kristen Miller, PhD

Federal Committee on Statistical Methodology Research and Policy Conference

Washington DC

October 25th, 2022

# Outline

- Background and context
  - COVID-19 pandemic
  - Open-text data: value and challenges
  - Item nonresponse detection: the technology and development of the model
- Evaluating the model: our approach
- Evaluation results
- Discussion/Next steps

# Background and context

# COVID-19 pandemic

- Numerous new COVID-19 related survey items

- Circumstances prevented our usual approach: in-depth cognitive interviewing to inform closed-ended online survey web probes

- Adapted and innovated our methods to include both closed and open-ended probes and experimental designs for post-hoc evaluations

# Open-text data: value and challenges

- Range of methodological uses for open-text data (Singer & Couper, 2017)

- Allows for responses without constraint (Schonlau & Couper, 2016) a particular advantage when little is known about a topic (Neuert et al., 2021, Scanlon, 2019; 2020)

- But higher response burden, more prone to item nonresponse, inadequate and irrelevant responses

- Coding and analysis can be labor intensive and time-consuming

- Recent advances in data science offer new efficiencies and opportunities

# Item nonresponse detection: prior work

- Categorizing item non-response
  - "nonproductive" responses (Behr et al., 2012)
  - Indirect (soft) versus direct (hard) refusals (Meitinger et al., 2021)

# Item nonresponse detection: prior work, cont'd

- Detecting item non-response

  - EvalAnswer* (Kaczmirek et al. (2017); available on GitHub)

    - **Complete non-response**: blank text box
    - **No useful answer**: "dfgjh"
    - **Don't knows**: "I have no idea"; "DK"; "I can't make up my mind"
    - **Refusals**: "no comment"; "see answer above"
    - **Other**: insufficient to code; "it depends"; "just do"; "just what it is"
    - **Single word**: "economy"
    - **Too fast**: < 2 seconds to answer

\* https://git.gesis.org/surveymethods/evalanswer

# Item nonresponse detection: prior work, cont'd

- Limitations of EvalAnswer

  - Relies on regular expressions (regex)

  - Missed some gibberish and don't know responses: "I dunno"; "no clue"

  - Flagged single word responses that are valid: "quarantine"; "furloughed"; "closings"

  - Flagged valid responses that include one of the rules:

    - "I have not bee unable to travel to see my grandsons who live away from me. I am **unsure** how this country is going to fare." [emphasis added]

  - Marked some non-response as valid:

    - "this is not a good question"; "I think my answer is self explanatory"

# Item nonresponse detection: Model development

- Trained a natural language processing (NLP) model to interpret responses.
  - Fine-tuned a Bidirectional Transformer for Language Understanding (BERT)* model using Simple Contrastive Sentence Embedding (SimCSE)**
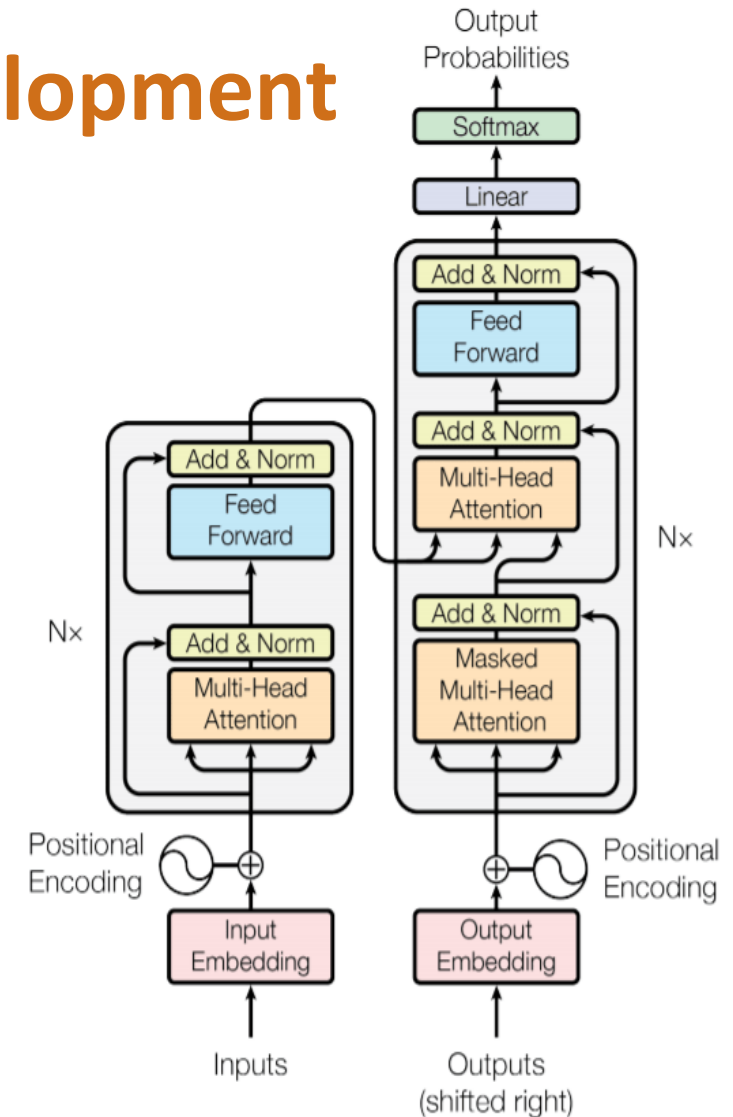- Refined training via human coding (active learning)



Figure 1: The Transformer - model architecture.

* https://arxiv.org/abs/1810.04805

** https://arxiv.org/abs/2104.08821

# Item nonresponse detection: Model development, cont'd

- Our working taxonomy:

  - **Complete non-response**: Blank text box

  - **Gibberish** or nonsensical: "dfgjh"

  - **Don't knows**: "I don't know"; DK; idk

  - **Refusals**: "no comment"; "Because"; "none"

  - **Other, high-risk**: non-useful response, non-codable

  - **Valid**: useful response, codable

- The model assigns a score (0-1) for the extent to which a response falls into each of the item non-response categories

# Model development: Active learning

- Round 1
  - 5 coders hand-coded 1,400 each, 200 overlapping with one other coder; full overlap for 500
  - Good consistency with most categories (gibberish, DKs, refusals)
  - Less consistency between valid versus "other, high risk" item nonresponse
  - Good results for identifying item nonresponse, but flagged more valids than we wanted
- Round 2:
  - 2 coders reviewed and arbitrated the results to retrain the model
  - Uncertainty retained in the model when warranted

# The data

- NCHS's Research and Development Survey (RANDS)
  https://www.cdc.gov/nchs/rands/index.htm

- RANDS During COVID-19 – Multi-round web/phone survey

- Topics: health, impacts of pandemic on health care access, COVID-19 related health care and behaviors

- Conducted using NORC at the University of Chicago's AmeriSpeak®, a probability-based panel representative of the US adult, English-speaking, non-institutionalized household population.

- Round 3 fielded May – June 2021: 5,458 Completes

  - 7,852 NORC's AmeriSpeak probability-based sample = 11.8% weighted cumulative response rate/69.5% completion rate

# Model evaluation: our approach

# Model evaluation: Phase 1

- Mixed-method evaluation of two web probe case studies

  - Quarantine probe (hand-coded data as source of truth)

  - Pandemic time reference probe (hand-reviewed sample as source of truth)

- Results (presented at AAPOR 2022)

  - Model did well at identifying "true" valids (high specificity); slightly less well identifying "true" item nonresponse (good sensitivity)

  - Outstanding issues:

    - Issue with false valid responses – "None", "none"

    - How well would the model perform on other non-COVID-19 related topics?

  - Subsequently, we retrained the model and carried out a Phase 2 evaluation

# Model evaluation: Phase 2

- Mixed-method evaluation of additional web probe case studies
  - Social distancing
  - Religion (new topic)

# Evaluation results

# Social distancing probe

- Social distancing survey questions:

  - In the last week, did you socially distance when you were…shopping, eating at a restaurant, etc. (total 7 randomized grid items)

  - [If yes, then] Did you do the following activities inside, outside, or both?

- Social distancing probe: When you were answering about social distancing in the previous questions, what were you thinking about?

- Full review of model-identified nonresponse (n=627); random sample (n=1,000) of valids

  - "Implied" sensitivity and specificity calculations

# Social distancing probe: evaluation results

| | Human-reviewed NR | Human-reviewed Valid | |
|---|---|---|---|
| **Model NR** | 450 | 177 | **627** |
| **Model Valid** | 100 = (25/1000)*3985 | 3885 = (975/1000)*3985 | **3985** |
| **Total** | 550 | 4062 | **4612** |

**Key take-away:**
**Model did a good job identifying "true" valids; slightly less well identifying "true" item nonresponse**

Sensitivity **82%** (450/550)

False valids (human-coded NR):
- "Recent activity"
- "EVERYTHING"
- "Being normal"
- "Don't do it as much"
- "Money"
- "I'm tired and I want to go to bed"

Specificity **96%** (3885/4062)

False NR (human-coded valid):
- "Safty" (and variations)
- "Save life"
- "lines in the market"
- "It is necessary but a pain."
- "Courtesy"
- "ITS COMMON CERDICY AND GO WITH THE THROW"

# Religion probe

- Religion survey question: Currently, how important is religion in your daily life? (very important, somewhat important, not important)

- Religion probe: Why do you say that?

- Full review of model-identified nonresponse (n=1,250); random sample (n=1,000) of valids

  - "Implied" sensitivity and specificity calculations

# Religion probe: evaluation results

| | Coded NR | Coded Valid | Total |
|---|---|---|---|
| **Model NR** | 298 | 952 | 1250 |
| **Model Valid** | 33 = (14/1000)*2350 | 2317 = (986/1000)*2350) | 2350 |
| **Total** | 331 | 3269 | 3600 |

**Key take-away:**
**Model did a good job identifying "true" item nonresponse; less well identifying "true" valids**

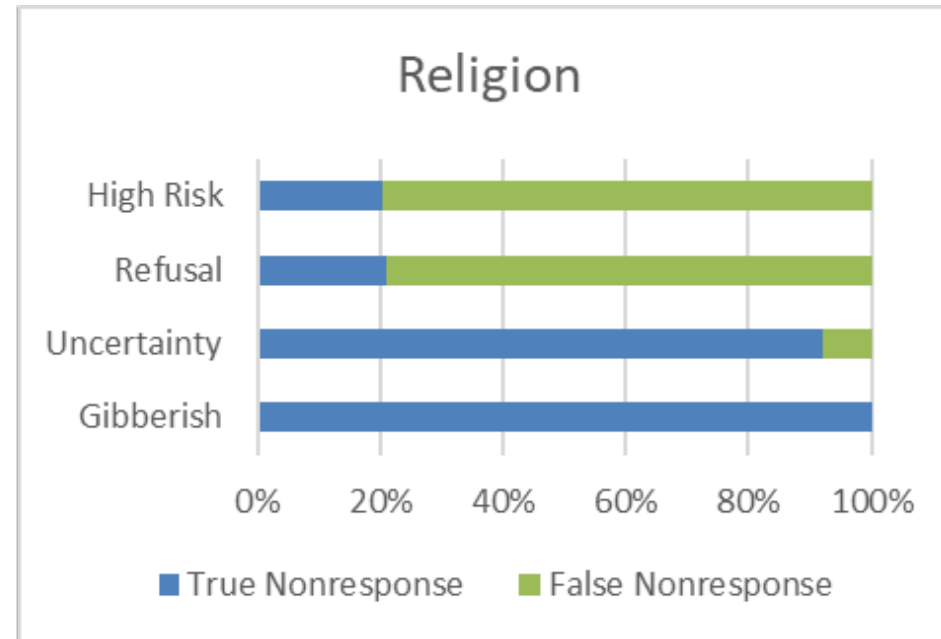Sensitivity **90%** (298/331)
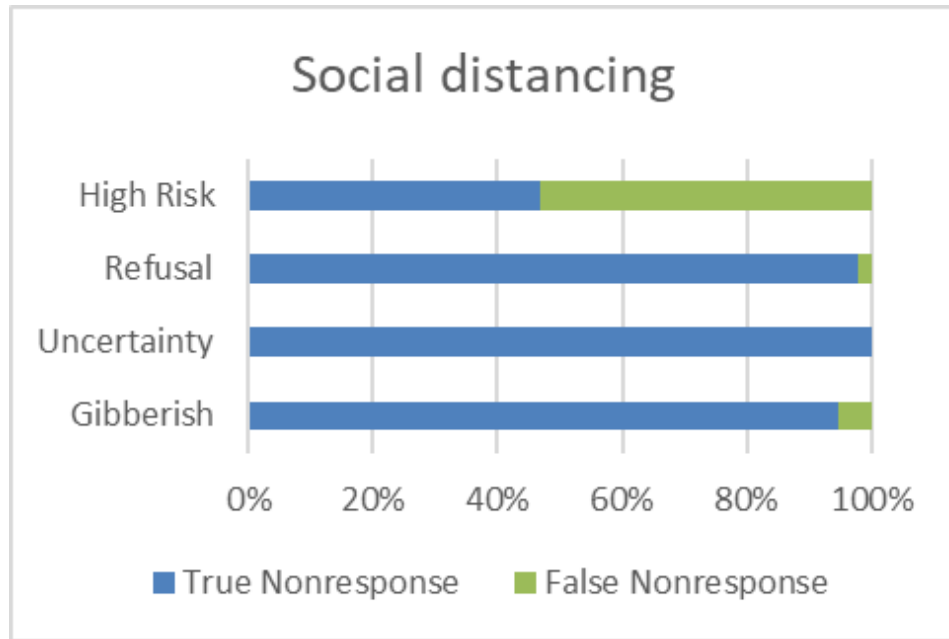
False valids (human-coded NR):
- "you asked me"
- "I JUST FEEL THAT WAY'
- "Guess"
- "Way of life"
- "Believe"

Specificity **71%** (2317/3269)

False NR (human-coded valid):
- "Faith"
- "It brings me peace"
- "I am not a religious person."
- "I worship mother earth. She is important"
- "My religion provides guides for living my life.  It encompasses my beliefs, goals and guidelines for living a good and right life."

# Distribution by type of item nonresponse



- Model error often concentrated in the High Risk category, as seen for Social distancing

- More error seen in Refusals for Religion

# Discussion/next steps

# Discussion/next steps

- Evaluation results show promise for our semi-automated item nonresponse detection model

- Next steps:

  - Release of a Semi-Automated Nonresponse Detector (SANDS) – a generalized model to share with others

  - Further evaluation and possible further training to understand and improve model performance on wider range of topics

  - Analysis to better understand the types and patterns of item nonresponse and possible subgroup differences

# Thank you!!

- Please contact us with any questions
  - Kristen Cibelli Hibben - kcibelli@cdc.gov
  - Zachary Smith – zsmith@cdc.gov
  - Travis Hoppe – thoppe@cdc.gov

**For more information contact:** Amanda Wilmot awilmot@cdc.gov

**Q-Bank:** providing access to survey question evaluation reports, question design and performance https://wwwn.cdc.gov/qbank/

**Q-Notes:** designed to facilitate the management and analysis of cognitive interviews https://www.cdc.gov/nchs/ccqder/products/qnotes.htm

**Centers for Disease Control and Prevention**

1600 Clifton Road NE, Atlanta, GA 30333

Telephone: 1-800-CDC-INFO (232-4636)/TTY: 1-888-232-6348

Visit: www.cdc.gov | Contact CDC at: 1-800-CDC-INFO or www.cdc.gov/info

# References

- Behr, D., Kaczmirek, L., Bandilla, W., & Braun, M. (2012). Asking probing questions in web surveys: which factors have an impact on the quality of responses? Social Science Computer Review, 30(4), 487-498.

- Kaczmirek, L., Meitinger, K., Behr., D. (2017). Higher data quality in web probing with EvalAnswer: a tool for identifying and reducing nonresponse in open-ended questions. (GESIS Papers, 2017/01). Köln: GESIS - Leibniz- Institut für Sozialwissenschaften.

- Schonlau, M. & Couper, M.P. (2016). Semi-automated categorization of open-ended questions. Survey Research Methods 10(2), pp. 143-152

- Singer, E. & Couper, M.P. (2017). Some methodological uses of responses to open questions and other verbatim comments in quantitative surveys. methods, data, analyses 11(2), pp. 115-134.

- Scanlon, P. J. (2019). The effects of embedding closed-ended cognitive probes in a web survey on survey response. Field Methods, 31(4), 328-343.

- Scanlon, P. (2020). Using targeted embedded probes to quantify cognitive interviewing findings. In P. C. Beatty, D. Collins, L. Kaye, J. Padilla, G. B. Willis & A. Wilmot (Eds.), Advances in questionnaire design, development, evaluation and testing, pp. 427–449.