

Identifying periods of careless responding in rating-scale surveys

Max Welz Andreas Alfons

ICORS, July 6, 2022

Ecolony

Erasmus University Rotterdam

Motivation

• Rating-scale surveys are ubiquitous in empirical research

- Not all respondents may respond accurately and truthfully
 - \rightarrow insufficient effort responding (IER; also called careless responding)

Careless responding

Definition: Careless Responding (Huang et al., 2012)

A response set in which the respondent answers a survey measure with low or little motivation to comply with survey instructions, correctly interpret item content, and provide accurate responses.

Identifying careless responding is important:

 \rightarrow Threat to internal validity (e.g. Huang et al., 2015)

 \rightarrow Already \leq 10% contamination is problematic (Arias et al., 2020)

Carelessness example: Random responding

Random responding: Tendency to randomly choose answer categories, regardless of item content

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
ICORS 2022 is awesome.					
I like chocolate.					
Oxygen is important.					
l am paid biweekly by leprechauns.				•	
Not getting bitten by a shark would be fun.					

L'IA

Careless responding as statistical outlier

- Careless respondents can be considered shape outliers example
- Existing outlier detection methods from robust statistics not suitable due to discrete nature of rating-scale data

- In long surveys, a large proportion of respondents may become careless towards the end due to fatigue (Bowling et al., 2021a)
 - $\longrightarrow\,$ Partial rowwise outlier
 - \rightarrow Goal: Identify <u>when</u> a given respondent responds carelessly (if at all)!

Proposed method (Part I)

Search for evidence of carelessness in two dimensions:

1 IER scores from autoencoder reconstruction (Kramer, 1992):

$$\mathsf{IER}_{ij} = \left(\frac{X_{ij} - \widehat{X}_{ij}}{L_j}\right)^2,$$

for reconstruction \widehat{X}_{ij} of response X_{ij} of *i*-th respondent to *j*-th item, for which there are L_j answer categories

Per-item response time (cf. Bowling et al., 2021b)

Illustration of change points



Proposed method (Part II)

Use Self-Normalization (Zhao et al., 2021) to detect change points in the two-dimensional series details

- Loops recursively over many nested segmentations of a series
- For each segmentation, calculate some test statistic
- If test statistic exceeds some threshold, there is a change point

- \longrightarrow Recursively narrow segment to isolate location of a change point
- \longrightarrow Asymptotic size control through the choice of the threshold

Simulation experiment

Generate 100 datasets details

• Sample 500 rating-scale observations for 15 constructs with 20 items each (random order)

 \longrightarrow Inspired by NEO PI-R (Costa and McCrae, 1992)

- Sample response times based on empirical times in Schroeders et al. (2022)
- Simulate carelessness onset from a Weibull distribution following Bowling et al. (2021a)
- Contaminate with various types of careless response styles

Simulated Carelessness onset for 40% contamination



Simulation results at level 2.5%



Using the Adjusted Rand Index (ARI; Hubert and Arabie, 1985) as performance measure, we find that our method

- accurately detects the change points for all response styles (ARI \geq 0.85)
- controls Type I error at levels 0.01, 0.025, 0.05

 \longrightarrow Results are fairly consistent across contamination levels (up to 100%)

Conclusions and outlook

- $\longrightarrow\,$ Proposed methodology seems promising for detecting careless periods
- \longrightarrow Next steps:
 - Tweaks in methodology: e.g., add third dimension based on longstring index to change point detection
 - Lots of simulations
 - Design a survey and collect emprical data

Conclusions and outlook

- \longrightarrow Robust methods for discrete data are underdeveloped
- $\longrightarrow\,$ Big interest in the social sciences for robustness
- \longrightarrow Potential for novel research ideas

My goal: develop theory and methods for robustness in discrete data

- \longrightarrow Measure transportation (e.g. Peyré and Cuturi, 2019) seems interesting for this purpose
- \longrightarrow I'd love to chat about with you about discrete robustness!

References

- Arias, V. B., Garrido, L., Jenaro, C., Martínez-Molina, A., and Arias, B. (2020). A little garbage in, lots of garbage out: Assessing the impact of careless responding in personality survey data. Behavior Research Methods, 52(6):2489–2505.
- Bowling, N. A., Gibson, A. M., Houpt, J. W., and Brower, C. K. (2021a). Will the questions ever end? person-level increases in careless responding during questionnaire completion. Organizational Research Methods, 24(4):718–738.
- Bowling, N. A., Huang, J. L., Brower, C. K., and Bragg, C. B. (2021b). The quick and the careless: The construct validity of page time as a measure of insufficient effort responding to surveys. Organizational Research Methods. In press.
- Costa, P. T. and McCrae, R. R. (1992). <u>Revised NEO personality inventory (NEO-PI-R) and NEO five-factor (NEO-FFI) inventory: prof</u> <u>Psychological Assessment Resources, Odessa, FL.</u>
- Goldfeld, K. and Wujciak-Jens, J. (2020). simstudy: Illuminating research methods through data generation. Journal of Open Source Software, 5(54):2763.
- Huang, J. L., Curran, P. G., Keeney, J., Poposki, E. M., and DeShon, R. P. (2012). Detecting and deterring insufficient effort responding to surveys. <u>Journal of Business and Psychology</u>, 27(1):99–114.

References (cont.)

- Huang, J. L., Liu, M., and Bowling, N. A. (2015). Insufficient effort responding: examining an insidious confound in survey data. Journal of Applied Psychology, 100(3):828–845.
- Hubert, L. and Arabie, P. (1985). Comparing partitions. <u>Journal of Classification</u>, 2(1):193–218.
- Kramer, M. A. (1992). Autoassociative neural networks. <u>Computers & Chemical Engineering</u>, 16(4):313–328.
- Peyré, G. and Cuturi, M. (2019). Computational optimal transport. <u>Foundations and Trends</u> in Machine Learning, 11(5-6):355–607.
- Scardapane, S., Comminiello, D., Hussain, A., and Uncini, A. (2017). Group sparse regularization for deep neural networks. Neurocomputing, 241:81–89.
- Schroeders, U., Schmidt, C., and Gnambs, T. (2022). Detecting careless responding in survey data using stochastic gradient boosting. <u>Educational and Psychological Measurement</u>, 82(1):29–56.
- Zhao, Z., Jiang, F., and Shao, X. (2021). Segmenting time series via self-normalization. <u>arXiv</u> preprint arXiv:2112.05331.

Ezafino

Appendix

Erafung 1/36

Example: Rating-scale data



Example: Random responding



Example: Straightlining



- Ordinal survey responses of n respondents to p items are collected in an n × p data matrix X, with possibly n < p
- The observations (respondents) are assumed to be i.i.d.
- We don't know when and where (if at all) IER occurs in $oldsymbol{X}$

Assumptions and main idea

Assumptions 1 and 2

- **1** The responses in **X** admit a low-dimensional representation, **S**, of dimension $q \ll p$. The value of q is known.
- 2 The survey that generated X is reliable in the sense that if there was no IER, X would accurately measure all constructs.

 \longrightarrow In the behavioral sciences, the number of constructs a survey measures is typically known

Idea: Reconstruct X from its low-dimensional representation S. Poor or near-perfect reconstruction might indicate IER (more on this later).

Auto-associative neural network (autoencoder)

Autoencoder (Kramer, 1992) to estimate **S** from **X**:

- Neural network with odd depth that attempts to reconstruct its input
- Central hidden layer is crucial: compresses the input to dimension q



Auto-associative neural network (autoencoder)

Network architecture:

- q nodes in bottleneck layer
- $1.5 \times p$ nodes in (de)mapping layers
- Activation functions:
 - \longrightarrow hyperbolic tangent in (de)mapping layers
 - \longrightarrow linear in bottleneck layer
- Run 100 epochs
- Use robust pseudo-Huber loss

Autoencoder: Regularization

We incorporate information on page-membership of items:

- \longrightarrow IER behavior can be expected to be similar within each page, but may differ between pages
- \rightarrow Use group-lasso regularization (cf. Scardapane et al., 2017)
 - Each group G_j holds items on the same survey page
 - Regularization is applied between input layer and mapping layer
- \rightarrow For instance, items from later pages (where IER can be expected to be higher) may be routed through different nodes in the mapping layer than items from earlier pages

Autoencoder: Regularization

For *m* pages and a parameter vector $\boldsymbol{\theta}$, the group-lasso penalty is

$$\Omega(\boldsymbol{ heta}) = \sum_{j=1}^m \sqrt{|G_j| \sum_{k \in G_j} \theta_k^2}$$



Autoencoder: Reconstruction error

Let \widehat{X}_{ij} be the reconstruction of response X_{ij} , i = 1, ..., n; j = 1, ..., p

 \longrightarrow The associated IER score is the reconstruction error

$$\mathsf{IER}_{ij} = \left(rac{X_{ij} - \widehat{X}_{ij}}{L_j}
ight)^2$$

with L_j denoting the number of answer categories of item j

Response times

Response times can be a good indicator of IER (Bowling et al., 2021b)...



Figure: Response times of a respondent who was instructed to respond accurately and truthfully to each item in a study of Schroeders et al. (2022)

Response times

... but sometimes they are not!



Figure: Response times of a respondent who was instructed to perform IER throughout a study of Schroeders et al. (2022)



Change point detection

Assumption 3

In periods in which a respondent engages in IER, there is a change point in either the IER scores or the response times, or both.

- \longrightarrow The IER scores and response times support or complement each other in identifying IER periods
- \longrightarrow Combine both quantities in a two-dimensional item series of length p and apply a multivariate method for change point detection

Estimation of multiple change points based on self-normalization (SNCP; Zhao et al., 2021)

- Let $\{\mathbf{Y}_t\}_{t=1}^p$ be a piecewise stationary series of dimension d
- Series has $m_o \ge 0$ (unknown) change points that partition it into $m_o + 1$ segments
- Each segment is piecewise stationary and obeys CDFs $F^{(i)}, i = 1, \dots, m_o + 1$
- CDFs are characterized by functionals $oldsymbol{ heta}_i = oldsymbol{ heta}\left(oldsymbol{F}^{(i)}
 ight) \in \mathbb{R}^d$
- \longrightarrow Breaks in $\{m{Y}_t\}_{t=1}^p$ are characterized by breaks in $m{ heta}$
- \longrightarrow Goal: Estimate number of change points, m_o , and their locations

For $1 \le a < b \le p$, put $\hat{\theta}_{a,b} = \theta\left(\hat{F}_{a,b}\right)$, where $\hat{F}_{a,b}$ is the empirical CDF of $\{\boldsymbol{Y}_t\}_{t=a}^b$. For a vector \boldsymbol{v} , denote its outer product by $\boldsymbol{v}^{\otimes 2} = \boldsymbol{v}\boldsymbol{v}^{\top}$. For $k \in \mathbb{N}$ such that $1 \le t_1 < k < t_2 \le p$:

We calculate the test statistic $T_p(t_1, k, t_2)$ based on a collection of nested windows covering k. Specifically, for $\epsilon \in (0, 0.5)$, define window size $h = \lfloor \epsilon p \rfloor$. For each $k = h, h + 1, \dots, p - h$, we define its corresponding nested window set $H_{1:p}(k)$ defined by

$$H_{1:p}(k) = \left\{ (t_1, t_2) \middle| t_1 = k - j_1 h + 1, \ j_1 = 1, \dots, \lfloor k/h \rfloor; \\ t_2 = k + j_2 h, \ j_2 = 1, \dots, \lfloor (p-k)/h \rfloor \right\}$$

For
$$k \in \{1, \ldots, p\}$$
 and $1 \leq s < e \leq p$, denote

$$W_{s,e} = \{(t_1, t_2) | s \le t_1 < t_2 \le e\}$$

and

$$H_{s:e}(k) = H_{1:p}(k) \cap W_{s,e}$$

and the subseries maximal test statistic by

$$T_{s,e}(k) = \max_{(t_1,t_2)\in H_{s:e}(k)} T_p(t_1,k,t_2)$$

 \rightarrow Algorithm 1 (next slide) uses these test statistics to estimate the number of change points, \widehat{m} , and their locations in $\{1, \ldots, p\}$

```
Algorithm 1: SNCP for multiple change point detection
```

```
Input: Time series \{\mathbf{Y}_t\}_{t=1}^p, threshold K_p, window size h = |p\epsilon|.
    Output: Estimated change-points set \hat{k} = (\hat{k}_1, \dots, \hat{k}_m)
    Procedure: SNCP(s, e, K_p, h), for 1 \le s < e \le p
    Initialization: SNCP(1, p, K_p, h)
   if e - s + 1 < 2h then
          Stop
 2
    else
 3
          \hat{k}^* = \arg \max_{k=s,\dots,e} T_{s,e}(k);
      if T_{s,e}(\hat{k}^*) < K_p then
 5
                Stop
 6
 7
          else
              \widehat{k} = \widehat{k} \cup \widehat{k}^*:
 8
              SNCP(s, \hat{k}^*, K_p, h);
 9
               SNCP(\hat{k}^*+1, e, K_p, h):
10
          end
11
12
    end
```

Algorithm 1 ...

- runs in $O(p/\epsilon^2)$ time
- identifies almost surely the correct number of change points and their correct locations, as $p \to \infty$

Under the null hypothesis of no change points, the SNCP test statistic $\max_{k=1,\dots,p} T_{1,p}(k)$ converges in distribution to a limit distribution $G_{\epsilon,d}$

- \rightarrow Offers (asymptotic) size control
- → For fixed ϵ , d, and $\alpha \in (0, 0.5)$, choose threshold K_p in Algorithm 1 as the (1α) -th quantile of $G_{\epsilon,d}$

Generating correlated ordinal variables

Sampling scheme of Goldfeld and Wujciak-Jens (2020):

- Goal: Sample independent X_j ∈ {1,..., K}, for j = 1,..., p, that jointly follow a positive semidefinite covariance matrix Σ
- \longrightarrow Ordinal variables with finite support (K answer categories)
 - Specify probabilities of choosing k-th answer category as

$$P_{j,k} = \mathbb{P}[X_j = k],$$

where $\sum_{k=1}^{K} P_{j,k} = 1$, for $j = 1, \ldots, p$

 \longrightarrow Implies distribution $X_j \sim F_j$, where, for $x \in \mathbb{R}$,

$$F_j(x) = \sum_{k=1}^{K} \mathbb{1}\{k \le x\} P_{j,k}$$

	Frahms
e	

Generating correlated ordinal variables

• For F_{logistic} the standard logistic CDF, put, for $k = 1, \dots, K$,

$$T_{j,k} = F_{\text{logistic}}^{-1}(F_j(k)),$$

with conventions $\mathcal{F}_{\mathsf{logistic}}^{-1}(0) = -\infty$ and $\mathcal{F}_{\mathsf{logistic}}^{-1}(1) = +\infty$

- Sample $(Y_1, \ldots, Y_p)^{\top} \sim \mathcal{N}_p(\mathbf{0}, \mathbf{\Sigma})$, then calculate probabilities $U_j = \Phi(Y_j)$ and quantiles $Z_j = F_{\text{logistic}}^{-1}(U_j)$, for $j = 1, \ldots, p$
- Finally, obtain the desired ordinal variables X_j as

$$X_j = 1 + \sum_{k=1}^{K} \mathbb{1}\{Z_j > T_{j,k}\}$$

 $\implies X_j \in \{1, \dots, K\}$ and $\mathsf{Var}[X_1, \dots, X_p] \approx \mathbf{\Sigma}$



Setup

Generation of 100 uncontaminated datasets:

- Generate n = 500 rating-scale observations
- 15 constructs with 20 items each ($\implies q = 15, p = 300$)
- $\bullet\,$ Constructs are mutually independent and items within the same construct have correlation of $\pm 0.7\,$
- All items use the following answer probabilities:

$$\frac{\mathbb{P}[X_j = 1]}{0.05} \quad \frac{\mathbb{P}[X_j = 2]}{0.25} \quad \frac{\mathbb{P}[X_j = 3]}{0.40} \quad \frac{\mathbb{P}[X_j = 4]}{0.25} \quad \frac{\mathbb{P}[X_j = 5]}{0.05}$$

- Items are in random order
- \rightarrow Inspired by the NEO PI-R (Costa and McCrae, 1992)

Simulation setup: IER onset

• Sample the item index at which IER onsets as i.i.d. draws from a Weibull distribution with location 240 (80% of all items), shape 2.2, and scale 20

 \rightarrow Inspired by Bowling et al. (2021a)

Setup: IER onset

 \longrightarrow 80% probability that IER starts before having answered 90% of all items (based on estimates in Bowling et al., 2021a)



Simulation setup: Response times

• Regular: Sample per-item response times (in seconds) as i.i.d. draws from a Weibull distribution with scale 6 and shape 2

 $\longrightarrow~$ Mean ≈ 5.3 and variance ≈ 2.8

 \rightarrow Emulates empirical observations in Schroeders et al. (2022)

- IER: Sample per-item response times (in seconds) as i.i.d. draws from a Weibull distribution with scale 2 and shape 1
 - $\longrightarrow\,$ Mean and variance of 2
 - \longrightarrow Inspired by Bowling et al. (2021b)

Setup: Response times

 \rightarrow Emulates empirical response times in Schroeders et al. (2022), also inspired by Bowling et al. (2021b)



Distribution of response times (in seconds)

Setup: Adding IER

- Vary the prevalence of IER: $\{0\%, 20\%, \dots, 100\%\}$
- Four types of IER: Random responding, straightlining, imperfect straightlining, pattern responding
- Each respondent who starts IER gets randomly assigned one of the four IER types

 \longrightarrow Each type is present in each contaminated dataset

Setup: Types of IER

- Random responding: Respond completely at random
- Straightlining: Respond always to the same, randomly determined answer category
- Imperfect straightlining: Straightlining, but with random variation of up to ± 1 answer category
- Pattern responding (based on Schroeders et al., 2022): Respond according to a fixed, randomly determined pattern

 \longrightarrow For example, 1-2-3-1-2-3, 4-5-4-5, ...

Evaluation measure

Adjusted Rand Index (ARI; Hubert and Arabie, 1985):

- ARI is a measure of classification performance
- Bounded between 0 (random classification) and 1 (perfect classification)
- \longrightarrow For each item, a simulated respondent is either IER or not
- \longrightarrow Calculate ARI of each respondent's series of per-item responses
- \longrightarrow Average the ARI over respondents

Results: ARI at level $\alpha = 0.025$



Results: ARI at level $\alpha = 0.025$



Results: ARI at level $\alpha = 0.025$



Simulation study: Size control

Results with size $\alpha = 0.01$:



Simulation study: Size control

Results with size $\alpha = 0.025$:



Simulation study: Size control

Results with size $\alpha = 0.05$:

