

From Individual Experience to Collective Evidence: An Incident-Based Framework for Identifying Systemic Discrimination

Jessica Dai, Paula Gradu, Inioluwa Deborah Raji, Benjamin Recht
University of California, Berkeley (alphabetical order)

Does some system that individuals interact with cause disproportionate harm to a meaningful subgroup?

- Individuals are described by covariates $X \in \mathbb{R}^d$, and groups G are defined by a function of covariates. Individuals can be in multiple groups at once; we know the base rate μ_G^0 for every group.
- Reports arrive sequentially. We *stop the algorithm* when we identify a group that has been harmed.

Given a reporting system, does any group report much more frequently than we would expect?

Example – vaccines (VAERS) and pharmaceuticals (FAERS)

- System: a particular vaccine or drug
- “Bad event”: medical adverse event (in general, or for a particular symptom)
- Reporting: reflects true “bad event” (not necessarily causal)
- Main question: “does this treatment disproportionately cause adverse events for a particular subgroup?”

Example – (algorithmic) decision-making

- System: a decision-making system, e.g. for loan allocations
- “Bad event”: an incorrect decision, e.g. denials to high-creditworthy individuals
- Reporting: correlated with (but not always equal to) “bad event”
- Main question: “does this decision-making system discriminate against a subgroup?”

Algorithm 1: Testing for overrepresentation in reports

Input: set of groups G with base preponderances $\{\mu_G^0\}_{G \in G}$, relative strength β , test level α , group size c

```
1 for  $t = 1, 2, \dots$  do
2   See an incident report  $X_t$ ;
3   for  $G \in G$  do
4     Update (log-)likelihood of harm for  $G$  depending on whether  $X_t \in G$ ;
5     Using differential privacy, choose  $\hat{G}^*$  to be the most impacted group based on reports seen thus far;
6     Test  $\hat{G}^*$  for harm at level  $\alpha$ ;
7     if  $\hat{G}^*$  is likely to be harmed then
8       Return  $\hat{G}^*$ .
```

Algorithm 2: Formal statement of Algorithm 1

Input: set of groups G with base preponderances $\{\mu_G^0\}_{G \in G}$, relative strength β , test level α , group size c

```
1 Initialize  $\omega_0^G = 0$  and  $\lambda^G = 1/2$  for all  $G \in G$ ;
2 Compute  $\tau = \ln(1 + \frac{1}{1-c})$  and  $\tilde{\alpha} = \max_{\gamma \in (0, \alpha)} (\alpha - \gamma) \exp(-\frac{1}{16} - \frac{1}{4} \sqrt{\ln(2/\gamma)})$ ;
3 for  $t = 1, 2, \dots$  do
4   for  $G \in G$  do
5     Update  $\omega_t^G \leftarrow \omega_{t-1}^G + \ln(1 + \lambda_t^G (1_{X_t \in G} - \beta \mu_G^0))$ ;
6     Sample  $\xi_t^G \sim \text{Lap}(4\tau\sqrt{2t})$ ;
7     Let  $z_t = \frac{1_{X_t \in G} - \beta \mu_G^0}{1 + \lambda_t^G (1_{X_t \in G} - \beta \mu_G^0)}$ ;
8     Let  $\lambda_{t+1}^G \leftarrow \text{Proj}_{[0, 1]}(\lambda_t^G + \frac{2}{2 - \ln(3)} \frac{z_t}{1 + \sum_{s \in [t]} z_s^2})$ ;
9   Let  $\hat{G}_t^* \leftarrow \arg \max_{G \in G} (\omega_t^G + \xi_t^G)$ ;
10  if  $\omega_t^{\hat{G}_t^*} \geq \ln(1/\tilde{\alpha})$  then
11    Return  $\hat{G}_t^*$ .
```

Main algorithmic idea: Sequential + multiple hypothesis testing for whether any group over-reports relative to their base rate.

Key questions:

- Correctness:** If we stop the test, is it true that our \hat{G}_t^* is actually reporting disproportionately often?
- Stopping time:** If there do exist groups that are harmed, how quickly can we identify them?

Some technical ingredients:

- “e-values”** (non-negative supermartingales) to ensure anytime-valid stopping (Alg. 2: L1, 5, 10)
- Differentially-private selection** to avoid Bonferroni over groups (Alg. 2: L2, 6, 9, 10)
- Online newton step** on λ_t^G to improve strength of the test (Alg. 2: L7, 8)

Detecting disproportionate reports of myocarditis in young men from COVID vaccines (VAERS data)

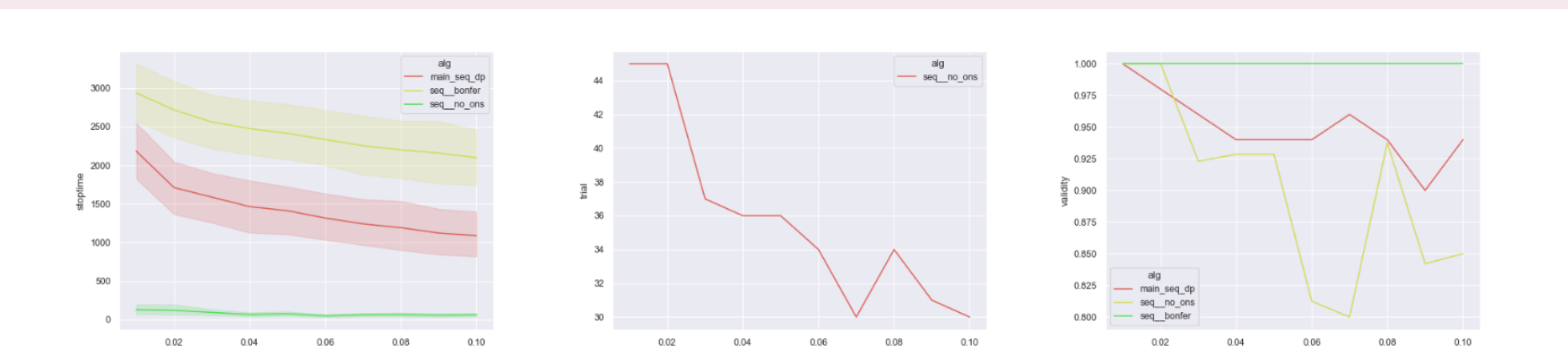


Figure 1: Results for incidents reporting Myocarditis or related symptoms after taking the COVID-19 vaccine over 50 trials. *Left*: Comparison of stop time across algorithms; *Center*: The non-ONS algorithm baseline fail rate again decreases with an increase in α ; *Right*: Comparison of validity across algorithms. $\alpha = [0.01 - 0.10]$, $\beta = 1.2$ for all.

Detecting disparity in loan denials to high-creditworthy individuals (HMDA data, simulated reports)

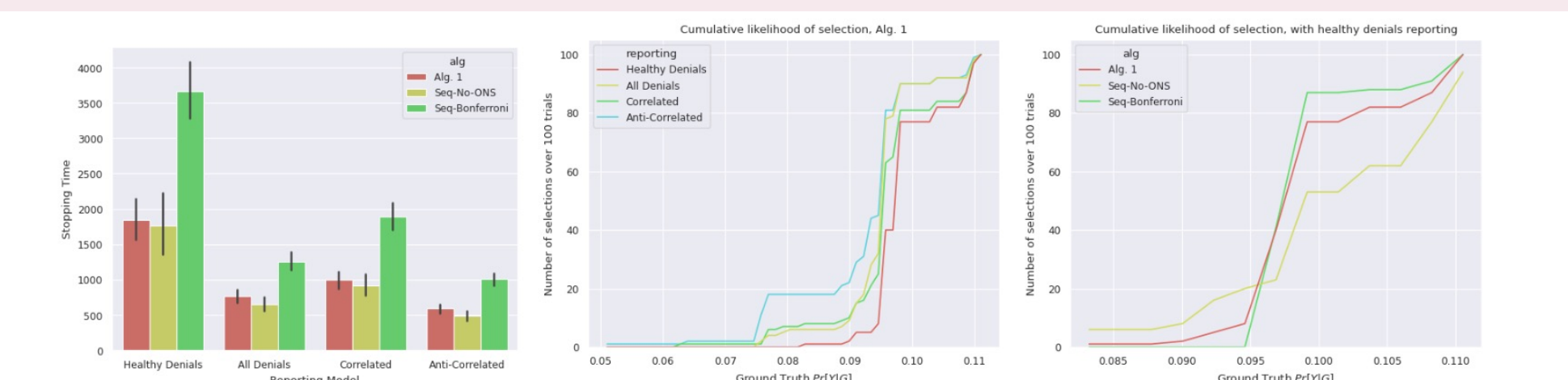


Figure 3: Impact of reporting models; 100 trials. *Left*: Stopping time by reporting model and algorithm. *Middle and right*: Ground-truth $\Pr[Y|G]$ vs cumulative selections. *Middle*: Alg. 2 across various reporting models; *Right*: comparison across algorithms for healthy denials reporting model. $\alpha = 0.01$, $\beta = 1.4$ for all.