Incentivizing Bandit Exploration: Recommendations as Instruments

Daniel Ngo, U of Minnesota Logan Stapleton, Vasilis Syrgkanis, Microsoft Research Zhiwei Steven Wu Carnegie Mellon

Example: Discriminatory lending

- Loan officer of a bank can offer low interest loan to a business client
- Unobserved client type correlates with both the baseline revenue and the treatment choice (which is biased against small businesses)

Bank wants to estimate the mean return (treatment effect) of these loans

loan officer's business selection type function



Model: Bank Loan Repeated Interaction

• For each round $t \in T$:

Recommendations sway some salespeople to offer low-interest loan to small businesses.

- New business comes in & assigned loan officer
- Bank recommends control ($z_t = 0$) or treatment ($z_t = 1$) to loan officer based on the history of returns on past loans
- Loan officer decides to offer loan ($x_t = 1$) or not ($x_t = 0$)
- Profit y_t is observed by the bank
- Estimate treatment effect $\hat{\theta}$ based on observed data

Bank's goal: Maximize net profit over all *T* rounds, estimate treatment effect

Loan officer's goal: Maximize their own expected profits

2/4

Recommendation as Instrument



Recommendations induce variability:

- Bank's recommendations incentivize loan officer to offer loan to small businesses
- Gather new information about loan returns

Instrument Variable Regression:

• Use *n* samples to estimate θ with $\hat{\theta}_n$:

Ш

$$\widehat{\theta}_n = \frac{\sum_{i=1}^n (y_i - \overline{y})(z_i - \overline{z})}{\sum_{i=1}^n (x_i - \overline{x})(z_i - \overline{z})}$$

3/4

Main Contributions

Multi-armed bandits algorithm with recommendations that incentivize exploration & act as instruments

Treatment effect estimation error bound after running our algorithm:



Regret bound for our algorithm: $R(T) = c + \tilde{O}(\sqrt{T})$

> prior-dependent constant

Incentivizing Bandit Exploration: Recommendations as Instruments

time horizon

Ш