Online Collaborative Prediction of Regional Vote Results

Vincent Etter, Emtiyaz Khan, Mattias Grossglauser, Patrick Thiran



DSAA — October 17, 2016 — Montréal, Canada

Data Opportunity

- Many countries adopt open government initiatives
- Several datasets published
 - Demographics
 - State affairs
 - Votes and elections
- Unique opportunity
 - Get a better understanding
 - Build **tools** useful to others

Voting Data

- News agencies, political parties, and polling institutes are all interested in understanding voting behaviors
 - Will the next vote pass easily?
 - What makes two regions vote similarly?
 - Where should we focus our efforts?

Dataset



- Vote results from Switzerland
 - Issue votes between 1981 and 2014
 - Outcome (% of "yes") at the municipality level

281 votes

•

• 13 features: voting recommendation of the main parties

• 2352 regions

• 25 features: languages spoken, demographics, etc.

Data available at http://vincent.etter.io/dsaa16

Similarities Between Results



Online Predictions

- On the day of the vote, regional results are released in **sequence**
 - Use **published results** to predict others

•

• ... and **refine** the prediction as more results are published?

Our Approach

- Use a matrix-factorization model to capture the **bi-clustering**
- Add region and vote features
 - Reduce the **cold-start** problem
 - More **interpretable**
- Build the model incrementally to assess the effect of each component

Our Model

$$y_{dn} = z_{dn} + \epsilon$$

$$z_{dn} = \mu_n + f_n(\boldsymbol{x}_d) + f_d(\boldsymbol{w}_n) + \boldsymbol{v}_d^T \boldsymbol{u}_n$$

$$\uparrow \qquad \uparrow \qquad \uparrow$$

$$f \qquad f \qquad f$$

$$f \qquad f$$

Our Models

Performance Evaluation

- Last 50 votes as test data
- Simulate 500 random reveal order
 - Last 10% of regions as test regions
 - **Observe** increasing number of regions
 - **Predict** result of test regions

Results



Bayesian VS Non-Bayesian



Final Model



Interpretation



Summary

- Individual models have different strengths
 - Vote features regression for cold start
 - Region features and bi-clustering when more observations
- Bayesian methods are useful
 - Proper hyperparameters setting
 - Accurate and interpretable results

Thank you!

Code and data available at

http://vincent.etter.io/dsaa16

Any questions?