# Diverse User Selection for Opinion Procurement

#### Yael Amsterdamer Oded Goldreich

Department of Computer Science







### Goals

- 1. Define user groups given high-dimensional data (on a scale)
- 2. Define "representativeness" measure of a user subset
  - Every user represents many groups (overlaps)
  - Many groups of different importance
- 3. Efficient computation of a good user subset
  - Minimizing manual effort
- 4. Optional customization

### **Previous work**

- Distance-based versus coverage-based diversity
- Diversity contexts:



#### **Our Solution – Brief Overview**

# **Profile Model**

User properties with scores in [0,1]



#### **Weighted Prioritization Framework**

Each group is assigned a **cover rate** and a **weight score** Natural options: uniform, proportional to size

Group cover scores are aggregated to a total score

## Grouping System

1-d clustering per property users are assigned to multiple groups

## **Explanations and Customizations**

**Explanations**: groups, selected users, group coverage

**Customization:** must(not) have groups, priority coverage, by weight manipulation

## **Max-Diversity Problem**

Maximize  $\sum_{G \in \mathcal{G}} \text{wei}(G) \cdot \min\{|U \cap G|, \text{cov}(G)\}$  for a selected subset U such that  $|U| \leq B$ 

Not penalizing over-representation, allowing under-representation

- NP-complete in B
- A greedy algorithm achieves a  $\left(1 \frac{1}{e}\right)$ -approximation in time  $O\left(B \cdot \max_{G \in \mathcal{G}} |G| \cdot \max_{u \in \mathcal{U}} \operatorname{groups}(u)\right)$

7

• Further practical optimizations via linked data structures

## **Explanations**

- Why user X was selected?
  - Top-k weighted groups represented by this user
- Importance of group Y?
  - Group name is indicative derived from profile data
  - Group weights



#### **Explanation-based Cutomization**

- Prioritize covering selected groups
  - Manually selected (1)
  - Ties are broken by other properties
- Target a better fit of a specific dimension
  - I.e., prioritize covering its groups
    - E.g., low, medium and high rating for Mexican food
- Focus on/filter out groups
  - Hard constraints

#### **Evaluation**

- Comparison to methods from IR, unsupervised learning
- Metrics:
  - Inherent quality of selected subset
    - Coverage of underlying (top-k) groups
    - Coverage of multi-dimensional meaningful groups, identified by mining
    - Similarity of distribution to source population
  - Diversity of task results
  - Client satisfaction

= 🔒 Podium	
Clara Taste buds satisfi	Selection Topic Coverage - 60% cantonese cuisine smoked duck chinese new year chef kong michelin star dim sum ritz carlton
What a great experience eating at this restaurant. The food was beautifully	service staff chinese tea char siew bamboo clams soup modern twist fine dining lobster meat abalo assistant manager noodles sea perch jasmine tea Random Selection Topic Coverage - 30%
<b>Do-yun</b> Best Cantonese K	chef kong (michelin star) (dim sum (ritz carlton) service staff) chinese tea (char siew) bamboo clams soup) (modern twist) (fine dining) (lobster meat) abalo (assistant manager) (noodles) sea perch (jasmine tea)

## **Experimental study**

#### User repositories: yelp, TripAdvisor



#### Intrinsic diversity (yelp)



#### Opinion diversity (yelp)



## Summary

- ✓ Weighted Framework for diverse user selection
- ✓ Population grouping system
- ✓ Diversity definition and efficient algorithm
- ✓ Explanations and customization
- Experimental studies indicate good balance between variance and coverage

# Thank you



