# Quality Assurance for Human Computation Based Recommendation

Master Defense Presentation

Michael Schwarz

June 13, 2017

Graz University of Technology

# Introduction

## Motivation

- Recommender systems are ubiquitous

## Motivation

- Recommender systems are ubiquitous
- These systems are usually based on knowledge

## Motivation

- Recommender systems are ubiquitous
- These systems are usually based on knowledge
- Reliable but expensive if entered by small number of experts

## Motivation

- Recommender systems are ubiquitous
- These systems are usually based on knowledge
- Reliable but expensive if entered by small number of experts
- Unreliable but cheap if entered by regular users

- Recommender systems are ubiquitous
- These systems are usually based on knowledge
- Reliable but expensive if entered by small number of experts
- Unreliable but cheap if entered by regular users
- Combine approaches to reliably and cheaply collect knowledge

- Design and implement a web-based generic recommender platform

- Design and implement a web-based generic recommender platform
- Add mechanisms to collect data from regular users

## Overview

- Design and implement a web-based generic recommender platform
- Add mechanisms to collect data from regular users
- Develop techniques to ensure the quality of the collected data

## Overview

- Design and implement a web-based generic recommender platform
- Add mechanisms to collect data from regular users
- Develop techniques to ensure the quality of the collected data
- Efficiently distribute tasks to users to improve the knowledge base

# Recommender Systems

## Recommender Systems

- There are different types of recommender systems

## Recommender Systems

- There are different types of recommender systems
- They all recommend products/items…

## Recommender Systems

- There are different types of recommender systems
- They all recommend products/items…
- …but use different techniques to find the best item(s)

## Recommender Systems

- There are different types of recommender systems
- They all recommend products/items…
- …but use different techniques to find the best item(s)
- Three types of systems are commonly used

- Collect information about the items (e.g., keywords)

- Collect information about the items (e.g., keywords)
- Find items similar to ones the user liked in the past

- Collect information about the items (e.g., keywords)
- Find items similar to ones the user liked in the past
- Idea: user preferences do not change

## Content-based Systems



- Collect information about the items (e.g., keywords)
- Find items similar to ones the user liked in the past
- Idea: user preferences do not change
- Advantage: independent of other users

- Collect information about the user

# Collaborative Filtering Systems

- Collect information about the user
- Find similar users

## Collaborative Filtering Systems



- Collect information about the user
- Find similar users
- Idea: people who liked the same things will like the same in the future
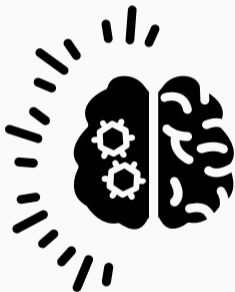
## Collaborative Filtering Systems



- Collect information about the user
- Find similar users
- Idea: people who liked the same things will like the same in the future
- Advantage: no understanding of the items necessary

- Explicit information about the items and user

- Explicit information about the items and user
- Find items that fulfill the user-given constraints

- Explicit information about the items and user
- Find items that fulfill the user-given constraints
- Idea: recommendation boils down to a constraint satisfaction problem

- Explicit information about the items and user
- Find items that fulfill the user-given constraints
- Idea: recommendation boils down to a constraint satisfaction problem
- Advantage: no history of the user is necessary

# A Generic Framework
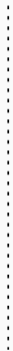
## Web-based Client-Server Model

- Subdivided into frontend and backend

## Web-based Client-Server Model

- Subdivided into frontend and backend
- Backend is based on the Spring Framework (Java)

## Web-based Client-Server Model

- Subdivided into frontend and backend
- Backend is based on the Spring Framework (Java)
- Frontend is mobile-friendly HTML5

## Web-based Client-Server Model

- Subdivided into frontend and backend
- Backend is based on the Spring Framework (Java)
- Frontend is mobile-friendly HTML5
- Parts are loosely coupled

Client

Server

Client

Server

| MessageHub (JavaScript) | ←→ | MessageHub (Java) |

Client

Server

| Message |
| --- |
| type (string) |
| content (JSON) |

Handler

MessageHub (JavaScript)

MessageHub (Java)

Handler

## Message Format
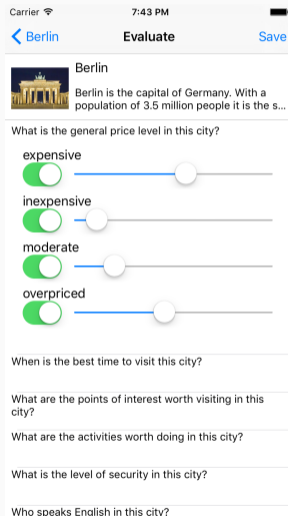
Message to register a new user

```
1 {
2   type: "register",
3   content : {
4     username: "michael",
5     password: "12345678",
6     email: "michael.schwarz@noreply.com"
7   }
8 }
```

## Multiple Frontends

- Loose coupling and easy API allows easy implementation of new frontends

- Loose coupling and easy API allows easy implementation of new frontends
- Bachelor Thesis: Implementation of a native iOS client
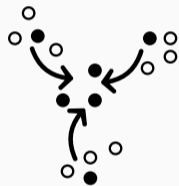
# Knowledge Acquisition

- Users do not like lengthy tasks

- Users do not like lengthy tasks
- Acquire knowledge from the user using small tasks (microtasks)

- Users do not like lengthy tasks
- Acquire knowledge from the user using small tasks (microtasks)
- Microtask has only one question

- Users do not like lengthy tasks
- Acquire knowledge from the user using small tasks (microtasks)
- Microtask has only one question
- 6 different types of microtasks

Item's support regarding one specific attribute

## Best matching item regarding one specific attribute



13

Best matching answer regarding one specific attribute

# Microtask #4

Weighted answers regarding one specific attribute

Implicit CAPTCHA

Binary decision

# Quality Assurance

- Users have to "earn" trust

- Users have to "earn" trust
- Score is influenced by CAPTCHAs, user behavior, etc.

- Users have to "earn" trust
- Score is influenced by CAPTCHAs, user behavior, etc.
- All contributions of the user are weighted with this score (0 % - 100 %)

- Users have to "earn" trust
- Score is influenced by CAPTCHAs, user behavior, etc.
- All contributions of the user are weighted with this score (0 % - 100 %)
- New or malicious users have minor to no influence on the knowledge base

- Depending on the human score, users get microtask with known answers (ground truth)

- Depending on the human score, users get microtask with known answers (ground truth)
- Similar to CAPTCHAs, but not seen as such by the user

- Depending on the human score, users get microtask with known answers (ground truth)
- Similar to CAPTCHAs, but not seen as such by the user
- Influence the human score (positively and negatively)

- Depending on the human score, users get microtask with known answers (ground truth)
- Similar to CAPTCHAs, but not seen as such by the user
- Influence the human score (positively and negatively)
- Classify an image, hard to do automatically

- We model the time it takes to answer a microtask

- We model the time it takes to answer a microtask
- Timings are matched using Kullback-Leibler distance

- We model the time it takes to answer a microtask
- Timings are matched using Kullback-Leibler distance
- Answers are weighted according to how well they fit

- We model the time it takes to answer a microtask
- Timings are matched using Kullback-Leibler distance
- Answers are weighted according to how well they fit
- Non-matching timings are discarded and decrease the human score

- Users can add new item, we have to cope with spam

- Users can add new item, we have to cope with spam
- CAPTCHAs only prevent automated spam

- Users can add new item, we have to cope with spam
- CAPTCHAs only prevent automated spam
- For a new item, we generate verification microtasks

- Users can add new item, we have to cope with spam
- CAPTCHAs only prevent automated spam
- For a new item, we generate verification microtasks
- If the community decides that an item does not belong to the recommender, it is removed

- We need knowledge for new items

- We need knowledge for new items
- Dynamic approach to calculate number of distributed microtasks

- We need knowledge for new items
- Dynamic approach to calculate number of distributed microtasks
- Loosely based on local working set algorithm for task scheduling

- We need knowledge for new items
- Dynamic approach to calculate number of distributed microtasks
- Loosely based on local working set algorithm for task scheduling
- Settle on minimum number of microtasks based on quality of the results

# Evaluation

- We conducted a <span style="color:orange">worldwide study</span>

- We conducted a worldwide study
- Users had to complete microtasks, evaluate items, and use the recommender

- We conducted a worldwide study
- Users had to complete microtasks, evaluate items, and use the recommender
- 1307 users (90.9 %) completed all tasks

- We conducted a worldwide study
- Users had to complete microtasks, evaluate items, and use the recommender
- 1307 users (90.9 %) completed all tasks
- Quality assurance led to recommendation improvement of >20 %

# Recommendation Quality Improvement

# Conclusion

- We developed a generic recommender framework for knowledge-based recommenders

# Summary

- We developed a generic recommender framework for knowledge-based recommenders
- We showed that users are willing to contribute through small tasks

# Summary

- We developed a generic recommender framework for knowledge-based recommenders
- We showed that users are willing to contribute through small tasks
- We presented automatic ways to ensure the quality of user content

**Human Computation Based Acquisition of Financial Service Advisory Practices**

*Alexander Felfernig, Michael Jeran, Martin Stettinger, Thomas Absenger, Thomas Gruber, Sarah Haas, Emanuel Kirchengast, Michael Schwarz, Lukas Skofitsch, Thomas Ulz*

FINREC'15

**Peopleviews: Human computation for constraint-based recommendation**

*Alexander Felfernig, Thomas Ulz, Sarah Haas, Michael Schwarz, Stefan Reiterer, Martin Stettinger*

ACM RecSys 2015 CrowdRec Workshop

**Human computation for constraint-based recommenders**

*Thomas Ulz, <u>Michael Schwarz</u>, Alexander Felfernig, Sarah Haas, Amal Shehadeh, Stefan Reiterer, Martin Stettinger*

Journal of Intelligent Information Systems 2016

**A Short Overview of the PeopleViews Mobile User Interface**

*Angela Promitzer, Alexander Felfernig, Michael Schwarz, Thomas Ulz, Amal Shehadeh, Sarah Haas*

TU Graz Technical Report 2016

**Thank you for your attention!**

THE BEST THESIS DEFENSE IS A GOOD THESIS OFFENSE

# Recommendation Quality Improvement without Ground Truth

## Human Score Calculation Example

| User | Human Score | Support | | | |
| --- | --- | --- | --- | --- | --- |
| | | Answer 1 | Answer 2 | Answer 1 (weighted) | Answer 2 (weighted) |
| User 1 | 1 | 0.8 | 0.3 | 0.8 | 0.3 |
| User 2 | 0.5 | 0.9 | 0.4 | 0.45 | 0.2 |
| User 3 | 0.5 | 0.6 | 0.5 | 0.3 | 0.25 |
| User 4 | 0 | 0.2 | 0.7 | 0 | 0 |
| *Sum* | **2** | 2.5 | 1.9 | 1.55 | 0.75 |
| *Average* | - | $\frac{2.5}{4} = 0.625$ | $\frac{1.9}{4} = 0.475$ | $\frac{1.55}{2} = 0.775$ | $\frac{0.75}{2} = 0.375$ |

**Table 1:** Four different users and their support values for Answer 1 and Answer 2.

## Optimal Number of Microtasks Example

|         | # of microtasks | Answered | Data is good | New # of microtasks |
|---------|-----------------|----------|--------------|---------------------|
| Cycle 1 | 10              | 4        | no           | $10 \times 1.5 = 15$ |

Goal: 5 answers

**Cycle 1** Start with 10 tasks $\rightarrow$ not enough, increase to 15

## Optimal Number of Microtasks Example

|          | # of microtasks | Answered | Data is good | New # of microtasks |
|----------|-----------------|----------|--------------|---------------------|
| Cycle 1  | 10              | 4        | no           | $10 \times 1.5 = 15$ |
| Cycle 2  | 15              | 11       | yes          | $15 \times 0.75 = 11$ |

Goal: 5 answers

**Cycle 1** Start with 10 tasks $\rightarrow$ not enough, increase to 15

**Cycle 2** 15 was enough, decrase to $15 \cdot 0.75 = 11$ tasks

## Optimal Number of Microtasks Example

|         | # of microtasks | Answered | Data is good | New # of microtasks |
|---------|:---------------:|:--------:|:------------:|:-------------------:|
| Cycle 1 | 10              | 4        | no           | $10 \times 1.5 = 15$ |
| Cycle 2 | 15              | 11       | yes          | $15 \times 0.75 = 11$ |
| Cycle 3 | 11              | 6        | yes          | $11 \times 0.75 = 8$ |

Goal: 5 answers

**Cycle 1** Start with 10 tasks $\rightarrow$ not enough, increase to 15

**Cycle 2** 15 was enough, decrease to $15 \cdot 0.75 = 11$ tasks

**Cycle 3** 11 was enough, decrease to $11 \cdot 0.75 = 6$ tasks

## Optimal Number of Microtasks Example

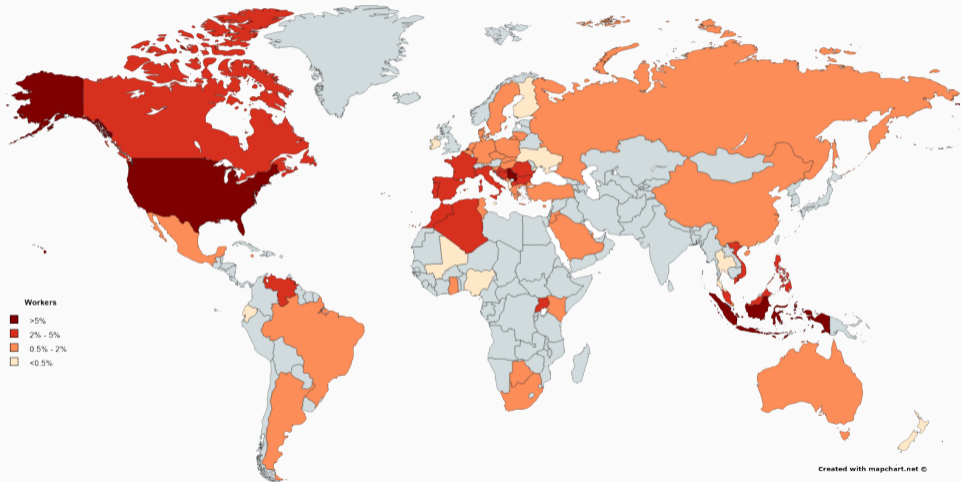|         | # of microtasks | Answered | Data is good | New # of microtasks |
|---------|-----------------|----------|--------------|---------------------|
| Cycle 1 | 10              | 4        | no           | $10 \times 1.5 = 15$ |
| Cycle 2 | 15              | 11       | yes          | $15 \times 0.75 = 11$ |
| Cycle 3 | 11              | 6        | yes          | $11 \times 0.75 = 8$ |
| Cycle 4 | 8               | 4        | no           | $8 \times 1.5 = 12$ |

Goal: 5 answers

**Cycle 1** Start with 10 tasks $\rightarrow$ not enough, increase to 15

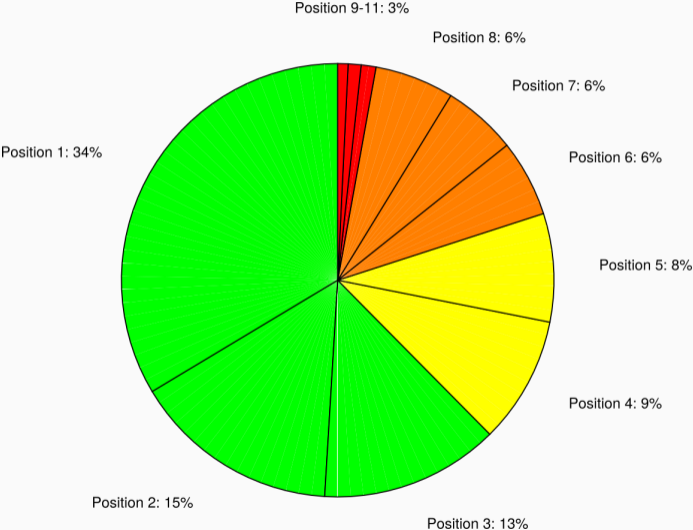**Cycle 2** 15 was enough, decrease to $15 \cdot 0.75 = 11$ tasks

**Cycle 3** 11 was enough, decrease to $11 \cdot 0.75 = 6$ tasks

**Cycle 4** 8 was not enough, increase to $8 \cdot 1.5 = 12$ tasks

# Worker Distribution



Workers
- >5%
- 2% - 5%
- 0.5% - 2%
- <0.5%

Created with mapchart.net ©

# Position of Chosen Item



Position 9-11: 3%

Position 8: 6%

Position 7: 6%

Position 6: 6%

Position 5: 8%

Position 4: 9%

Position 3: 13%

Position 2: 15%

Position 1: 34%

# Recommendation Screen

Refine by

## "City"

← Back

⚡ Play against the community

| Item | | | Compare |
|---|---|---|---|
| | **New York City** The City of New York, often called New York City or simply New York, is the most populous city in the United States.[1] ... | 🏆 Best in 2 attributes... | ☐ |
| | **Paris** Paris, the cosmopolitan capital of France, is one of the largest agglomerations in Europe, with 2.2 million people livin... | 🏆 Best in Museums | ☐ |
| | **Sydney** Sydney is the state capital of New South Wales and the most populous city in Australia and Oceania. Located on Australia... | | ☐ |
| | **Berlin** Berlin is the capital of Germany. With a population of 3.5 million people it is the second most populous city and the se... | | ☐ |
| | **Moscow** Moscow is the capital and the largest city of Russia. It is a major political, economic, cultural, and scientific center... | | ☐ |
| | **Mumbai** Mumbai (also known as Bombay, the official name until 1995) is the capital city of the Indian state of Maharashtra. Mumb... | | ☐ |
| | **Beijing** Beijing is the capital of the People's Republic of China and one of the most populous cities in the world. | | ☐ |
| | **Tokyo** Tokyo, officially Tokyo Metropolis, is one of the 47 prefectures of Japan, and is both the capital and largest city of J... | 🏆 Worst in Everyone | ☐ |

« | 1 - 8 | 9 - 11 | »

🗐 Compare items   💾 Save filter