

jobsta

*Study - Job - Money*

Project for  
Introduction to Semantic Systems  
(188.399-2019W)

Group 01

Cem Bicer (01425692)  
Helmuth Breitenfellner (08725866)  
Laszlo Kiraly (09227679)  
Gerald Weber (00125536)

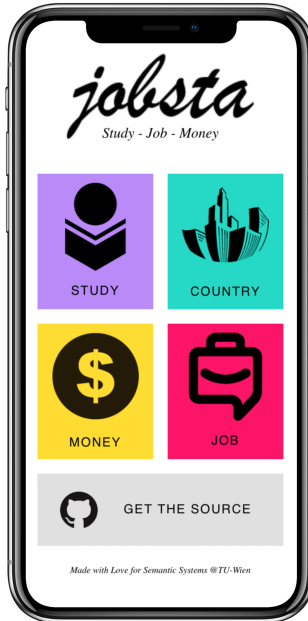
2020-01-31

The logo for 'jobsta' is written in a black, cursive, handwritten-style font. The letters are connected, with a prominent 'j' and a stylized 'a'.

*Study - Job - Money*

- For Software Developers and Data Scientists
- Asks for experience, age, location
- Answers to following questions:
  - *What shall I study?*
  - *Where shall I work?*
  - *What shall I practise?*
  - *How can I improve?*

# The Mobile App



## Data Sources

- **Kaggle User Survey**  
Data Scientists, Country, Job Role, Programming Language, Income
- **StackOverflow User Survey**  
Software Developer, Country, Job Role, Programming Language, Income
- **GitHub Repositories Data**  
Repository URL, Popularity, Programming Language, Issues
- **TISS Lectures**  
Lectures, Lecturer, Description, Programming Language

## Kaggle Survey

- <https://www.kaggle.com/c/kaggle-survey-2019>
- Used Jupyter Notebook for Pre-Processing
- Created RDF directly from Python
- Private data was not usable for our case
- **Challenge:** harmonize data

## StackOverflow Survey

- <https://insights.stackoverflow.com/survey/2018>
- Used Python for Pre-Processing
- Created RDF directly from Python
- **Challenge:** Excel, Numbers and TextMate could all not open the csv file (>90.000 entries) properly

## GitHub Repositories Data

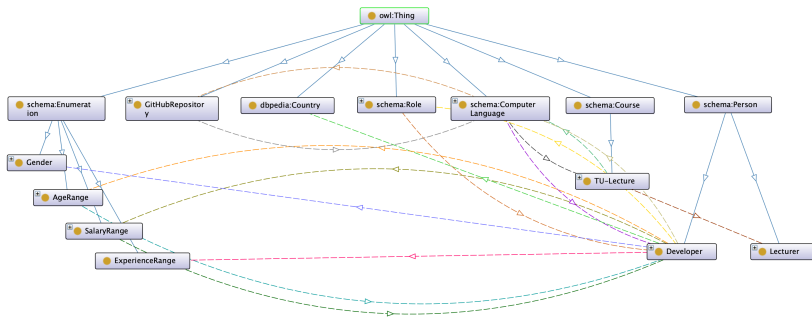
- <http://ghtorrent.org/>
- Used Bash & R Script for Pre-Processing
- Created RDF directly from Python
- **Challenge:** *huge* data archive (>100GB) had to be filtered / preprocessed

## TISS Lectures

- <https://tiss.tuwien.ac.at/course/courseList.xhtml?dswid=6403&dsrid=238>
- Used Python Script
- Created RDF directly from Python (using `rdflib`)
- **Challenge:** web scraping, identifying the programming language from text



# Ontology #1



## Ontology #2

- Created with Protégé
- Reusing existing Ontologies
  - [schema.org](http://schema.org)
  - [dbpedia.org](http://dbpedia.org)
- Entites:
  - Developer
  - dbpedia:Country
  - schema:Course
  - GitHubRepository
  - schema:ComputerLanguage
  - ...
- Object properties:
  - dealsWith
  - developsIn
  - schema:homeLocation
  - ...

# Harmonize Data I

- Age Ranges
  - Different Age Ranges
- Salary vs. Salary Range
  - Salary Range in Kaggle
  - Salary Value in Stackoverflow
- Roles
  - Combined from Surveys into List
  - e.g. Frontend Developer -> Software Engineer
  - ... C-Suite Executive -> Manager

## Harmonize Data II

- Countries
  - dbpedia linked to external data
- Gender
  - Single Selection in Kaggle
  - Multiple Selections in Stackoverflow
- Computer Language
  - Combined from Surveys into List
  - Field in Github Repository
  - Extracted from TISS Lecture Description

## SPARQL Queries #1

“As a developer I live in (Austria) and I want more than (150000 USD per year). What courses at TU Wien deal with programming languages which high-earners are using?”

# SPARQL Queries #1 (cont.)

```
1 ▼ PREFIX group1: <http://www.semanticweb.org/sws/ws2019/group1#>
2 PREFIX schema: <http://schema.org/>
3 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
4 PREFIX dbpedia: <http://dbpedia.org/resource/>
5
6 SELECT DISTINCT ?lecture ?language
7 ▼ WHERE {
8     ?lecture group1:dealsWith ?language .
9     ?developer group1:developsIn ?language .
10    ?salaryRange group1:maxSalary ?salary .
11    ?developer group1:hasSalaryRange ?salaryRange .
12    ?developer schema:homeLocation dbpedia:Austria .
13    FILTER (?salary > "150000"^^xsd:integer)
14 }
15 LIMIT 25
```

## SPARQL Queries #2

“As a (Data Scientist) I live in (Austria) and I can program in (Java) and I want more than (62000 USD per year). Should I stay or should I go?”

## SPARQL Queries #2 (cont.)

```
11 ASK
12 v WHERE {
13     ?developer a group1:Developer .
14     ?developer schema:homeLocation ?country .
15     ?developer group1:developsIn ?language .
16     ?developer group1:hasRole ?role .
17 v {
18     SELECT ?country (AVG(?avgRange) as ?averageK)
19 v WHERE {
20         ?developer a group1:Developer .
21         ?developer group1:hasRole ?role .
22         ?developer schema:homeLocation ?country .
23         ?developer group1:developsIn ?language .
24         ?developer group1:hasSalaryRange ?salaryRange .
25         ?salaryRange group1:minSalary ?minSalary .
26         ?salaryRange group1:maxSalary ?maxSalary .
27         BIND ((?minSalary + ?maxSalary)/2 AS ?avgRange)
28     }
29     GROUP BY ?country
30 }
31 v {
32     SELECT ?country (AVG(?salaryValue) as ?averageS)
33 v WHERE {
34         ?developer a group1:Developer .
35         ?developer group1:hasRole ?role .
36         ?developer schema:homeLocation ?country .
37         ?developer group1:developsIn ?language .
38         ?developer group1:salary ?salaryValue .
39     }
40     GROUP BY ?country
41 }
42 BIND ((?averageK + ?averageS)/2 as ?average)
43 FILTER (?language = group1:Java && ?country = dbpedia:Austria && ?average > "62000"^^xsd:integer && ?role = group1:Data_Scientist)
44 }
45 GROUP BY ?country ?average
```



## Lessons Learned

- Iterative process to come up with final idea
- Scraping TISS: no ID access to fields
- <http://schema.org> not equal to <http://www.schema.org>
- GraphQL Github API vs. Database Dump
- Harmonizing data can be tedious

Questions?

Thank you for your attention!