

Information Retrieval  
Academic Year 2022-23  
**Course Project**  
Prof. Monica Landoni

For this assignment you will work in pairs to implement a working prototype of an information retrieval system for a specific task and user needs. The goal of this assignment is to apply the concepts and tools you have learned in the class to a practical, real world, application.

### **Description**

Each group was assigned a topic, together with an example website you need to crawl your data from. You need to build a system that gathers a large collection of samples associated with the topic and enable the search over this collection. To build the collections, you need to crawl multiple sources (websites), similar in content and topic to the website suggested. The system must provide an interface for searching, browsing, and presentation of the data to the user.

The system should also have two additional features. Additional features are described in the document on iCorsi and categorized into two groups: complex and simple features. You can choose which features you want to implement yourself, provided there is at least one feature from each of the groups (two in total).

### **Submission procedure and evaluation**

An important step is the evaluation of your tool. Three students of the class, of your choice, will act as test users of your system and help you in evaluating your system from the point of view of the user experience. System evaluation (i.e. recall and precision) is not required. Please, coordinate with the other groups in order to meet the deadline.

You have to produce a report (max. 10 pages, including cover) of your work and of the evaluation. It will contain a concise explanation of how you tackled the design and implementation of the system. The code of the project and the report need to be submitted via iCorsi by end of day 7<sup>th</sup> December.

# Projects

N.	Project	Guidelines	Example Website
1	<b>IKEA hacks</b>	Choose at least 3 different websites with IKEA hacks \ articles about IKEA furniture \ reviews of IKEA products.	<a href="https://ikeahackers.net/category/hacks">https://ikeahackers.net/category/hacks</a>
2	<b>TV Series Episodes</b>	Choose at least 3 different websites about TV series'	<a href="https://www.tvmaze.com/">https://www.tvmaze.com/</a>
3	<b>Volunteer camps</b>	Choose at least 3 websites listing volunteer camps or opportunities.	<a href="https://www.workcamps.sci.ngo/icamps/welcome.html">https://www.workcamps.sci.ngo/icamps/welcome.html</a> or <a href="https://www.volunteerworld.com/en">https://www.volunteerworld.com/en</a>
4	<b>Hairdressers</b>	Choose at least 3 websites and then focus on a specific city, any other kind of professional (plumbers, lawyers etc) is also ok. If you want to mix and match (like plumbers, electricians and painters, everything for the home, or hairdressers + aestheticians) it's also fine.	<a href="https://www.fresha.com/lp/en/bt/hair-salons/in/gb-london">https://www.fresha.com/lp/en/bt/hair-salons/in/gb-london</a>
5	<b>Clothing</b>	Choose at least 3 websites listing clothing (Asos, Zalando, Shein...)	<a href="https://www.asos.com/">https://www.asos.com/</a>
6	<b>Short-term rentals</b>	Choose at least 3 websites, then focus on a specific city\area.	<a href="https://www.airbnb.co.uk/">https://www.airbnb.co.uk/</a>
7	<b>Textbooks</b>	Choose at least 3 websites listing textbooks (any subject)	<a href="https://open.umn.edu/opentextbooks">https://open.umn.edu/opentextbooks</a>
8	<b>Charities</b>	Choose at least 3 websites listing charities.	<a href="https://www.guidestar.org/NonprofitDirectory.aspx">https://www.guidestar.org/NonprofitDirectory.aspx</a>
9	<b>Indie Games</b>	Choose at least 3 websites listing indie videogames (commercial platforms like Steam are ok but only crawl indie games)	<a href="https://itch.io/games">https://itch.io/games</a>
10	<b>Stocks</b>	Choose at least 3 different stock exchanges.	<a href="https://www.borsaitaliana.it/borsa/azioni/listino-a-z.html?initial=A&amp;lang=en">https://www.borsaitaliana.it/borsa/azioni/listino-a-z.html?initial=A&amp;lang=en</a>
11	<b>Travel experiences</b>	Choose at least 3 websites listing travel experiences/tour, then focus on a specific city\area\country.	<a href="https://www.getyourguide.co.uk/london-157/">https://www.getyourguide.co.uk/london-157/</a>
12	<b>PhD positions</b>	Choose at least 3 websites listing PhD positions, any subject.	<a href="https://academicpositions.com/jobs/position/phd">https://academicpositions.com/jobs/position/phd</a>
13	<b>Content creators</b>	Choose at least 3 websites for Content Creators (like Patreon, Ko-fi...but please keep it safe for work!)	<a href="https://www.patreon.com/en-GB">https://www.patreon.com/en-GB</a>
14	<b>Concerts</b>	Choose at least 3 websites, then focus on a specific country, or 3	<a href="https://www.ticketmaster.ch/music/concerts/1221/events?language=en-us">https://www.ticketmaster.ch/music/concerts/1221/events?language=en-us</a>

		smaller countries (Switzerland is a small country!)	
15	<b>Pokemon</b>	Choose at least 3 websites listing pokemons.	<a href="https://www.pokemon.com/us/pokedex/">https://www.pokemon.com/us/pokedex/</a>
16	<b>Pets for adoption</b>	Choose at least 3 websites listing pets for adoption.	<a href="https://www.petfinder.com/">https://www.petfinder.com/</a>
17	<b>Nannies/babysitters</b>	Choose at least 3 websites listing nannies\babysitters, then focus on a specific city\areas.	<a href="https://www.childcare.co.uk/find/Babysitters">https://www.childcare.co.uk/find/Babysitters</a>
18	<b>Motorcycles</b>	Choose at least 3 websites listing motorcycles.	<a href="https://www.motorcycle.com/">https://www.motorcycle.com/</a>
19	<b>Mangas</b>	Choose at least 3 websites listing Mangas. Manwha or Manhwa (Korean or Chinese comics) are also ok!	<a href="https://mangatx.com/manga/">https://mangatx.com/manga/</a>
20	<b>Second-hand items</b>	Choose at least 3 websites offering second-hand items.	<a href="https://www.secondhand.org.uk/">https://www.secondhand.org.uk/</a>
21	<b>Football teams</b>	Choose at least 3 websites listing football teams, not necessarily from the same country.	<a href="https://footballdatabase.com/">https://footballdatabase.com/</a>
22	<b>Medications</b>	Choose at least 3 websites listing medications.	<a href="https://www.drugs.com/drug_information.html">https://www.drugs.com/drug_information.html</a>
23	<b>Knitting patterns</b>	Choose at least 3 websites listing knitting\crochet patterns.	<a href="https://www.garnstudio.com/search.php?action=browse&amp;c=home&amp;lang=en">https://www.garnstudio.com/search.php?action=browse&amp;c=home&amp;lang=en</a>

# Groups

	Partner 1	Partner 2	Project number
1	Agostino Monti	Marco Farace	2
2	Albert Cerfeda	Alessandro Gobbetti	13
3	Alen Sugimoto	Mattia Monari	7
4	Alessandro Cagnani	Diell Kryeziu	5
5	Alessandro Cravioglio	Alfio Vavassori	19
6	Andrea Cardia	Fabrizio De Castelli	16
7	Andrea Gualandris	Alan Copa	6
8	Andrea Prato	Samuel Corecco	14
9	Anton Tanev	Manuele Jelmini	1
10	Alessio Cordivani	Roberto Palmieri	4
11	Bojan Lazarevski	Daniela Gjorgjieva	20
12	Catarina Morais	Michele Zucchi	8
13	Enrico Di Pietro	Giacomo Solaro	17
14	Francesco Casarella	Marco Biasion	15
15	Francesco Costa	Arnaud Fauconnet	21
16	Johan Jacob	Dylan Reid Ramelli	11
17	Kelvin Likollari	Gerald Prendi	18
18	Kevin Spahiu	Fabio Zampielo	10
19	Mak Fazlic	Harkeerat Sawhney	3
29	Ramazan Tafa	Daniel Klimas	9
21	Razvan Petrica Onciu	Elisa Spinelli	12
22	Zhang Xufeng	Arash Mowdoudi	23
23	Bahram Ismayilov	Hoormazd Pirayeshfar	22

## Simple Features

**Results presentation:** results should be presented in a tabular format, so that many results could be seen at the same time. Each table cell should contain appropriate information for your project.

**Filtering:** in addition to being able to search by title, an user should be able to filter the results based on at least 3 relevant attributes for your project.

**Results Snippets:** present result snippets of each retrieved result (maximum 2-3 lines) in a kind of “Google style”, with query terms highlighted.

## Complex Features

**Automatic Recommendation:** In addition to the relevant search results pertaining to the user query, the user should also be suggested “similar” products based on, say, category, description, price etc. The ordering among the recommended items is not important. However, you should mention how did you arrive at the recommendations. For this purpose, you can use any open-source recommenders available.

**User Relevance Feedback:** after presenting the search results to an user, the user may provide a positive or negative feedback on the results (i.e. mark relevant and irrelevant documents). Based on this feedback the search results have to be updated and presented again.

**Results clustering:** the system should group results into topics. There should be a possibility for expanding/shrinking a topic to show all the results related to the topic. In addition, topics should be sorted in a descending order by the number of relevant results under them.