



SOCIAL MEDIA DATA MINING AND USER EMPOWERMENT

Ali Padyab, PhD
Information Systems

Department of Computer Science, Electrical and Space Engineering



AGENDA

- Introduction of myself
- Social media data mining
- Why do we need privacy? The role of awareness
- Introduction to USEMP project tools: DataBait (Collaboration between LTU and CEA)
- Demonstration of inferences in Facebook
- Implications for digitalization
- Q&A



A large, detailed iceberg floating in a dark blue sea, with its reflection clearly visible on the water's surface. The sky is a pale, hazy blue.

INTRODUCTION

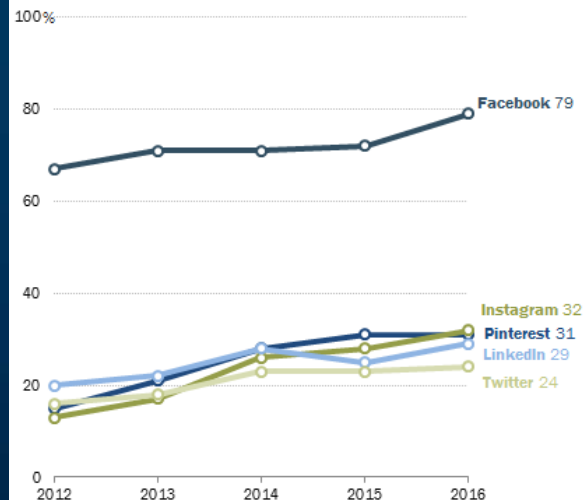
- Ali Padyab, PhD
- Field of research: Information Privacy, end user attitude, Privacy Enhancing Tools, IoT, Smart city



SOCIAL MEDIA USERS

Facebook remains the most popular social media platform

% of online adults who use ...

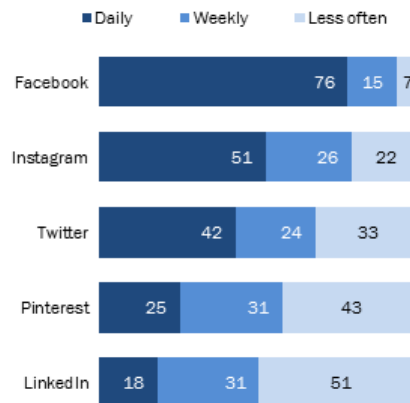


Note: 86% of Americans are currently internet users
 Source: Survey conducted March 7-April 4, 2016.
 Social Media Update 2016

PEW RESEARCH CENTER

Three-quarters of Facebook users and half of Instagram users use each site daily

Among the users of each social networking site, % who use these sites ...



Note: Do not know/refused responses not shown.
 Source: Survey conducted March 7-April 4, 2016.
 Social Media Update 2016

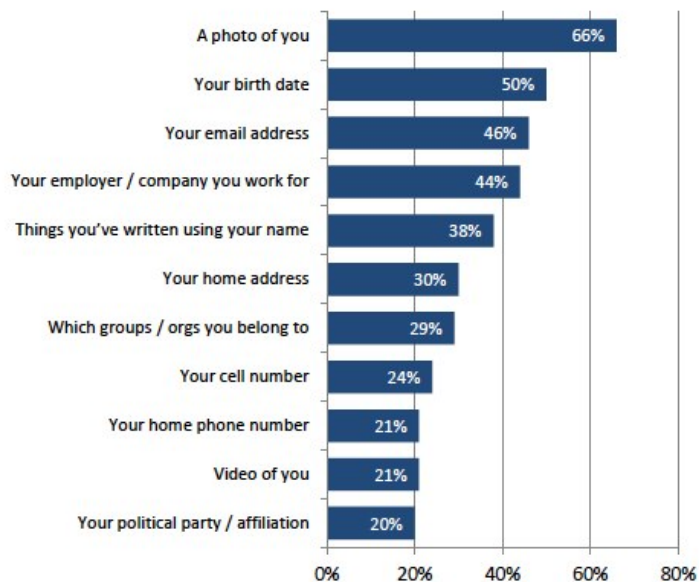
PEW RESEARCH CENTER



PERSONAL INFORMATION ONLINE

Personal information online

% of adult internet users who say this information about them is available online

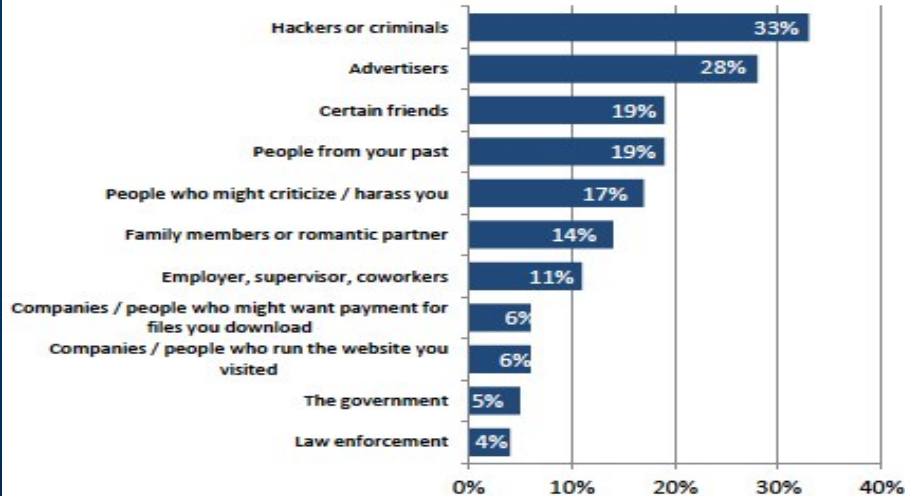


Source: Pew Research Center's Internet & American Life Project Omnibus Survey, conducted July 11-14, 2013, on landline and cell phones. N=792 for internet users and smartphone owners. Interviews were conducted in English on landline and cell phones. The margin of error on the sample is +/- 3.8 percentage points.



Who users try to avoid

% of adult internet users who say they have used the internet in ways to avoid being observed or seen by ...



Source: Pew Research Center's Internet & American Life Project Omnibus Survey, conducted July 11-14, 2013, on landline and cell phones. N=792 for internet users and smartphone owners. Interviews were conducted in English on landline and cell phones. The margin of error on the sample is +/- 3.8 percentage points.



ACCESS TO OSN DATA

Broad concern about government and third-party access to info on social networking sites

% Among adults ages 18 and older who use social networking sites

How concerned are you that some of the info you share on social networking sites might be accessed by ___ without your knowledge?

	The government	Third parties like advertisers or businesses
Very concerned	37	35
Somewhat concerned	34	45
Not too concerned	25	17
Not at all concerned	4	2

Source: Pew Research Privacy Panel Survey, January 2014. N=607 adults, ages 18 and older.

PEW RESEARCH CENTER



WHAT INFORMATION COLLECTED ACCORDING TO PRIVACY POLICIES

facebook

What kinds of information do we **collect**?

- Things you do and information you provide.
- Things others do and information they provide.
- Your networks and connections.
- Information about payments.
- Device information.
- Information from websites and apps that use our Services.
- Information from third-party partners.
- Facebook companies

Google

What kinds of information do we **collect**?

- Information you give us
- Information we get from your use of our services
 - Device information
 - Log information
 - Location information
 - Unique application numbers
 - Local storage
 - Cookies and similar technologies

WHAT ARE USES ACCORDING TO PRIVACY POLICIES?

facebook.

Google

How do we **use** this information?

- Provide, improve and develop Services.
- Communicate with you.
- Show and measure ads and services.
- Promote safety and security.

How we **use** information we collect?

- ...to provide, maintain, protect and improve them, to develop new ones, and to protect Google and our users. ... to offer you tailored content – like giving you more relevant search results and ads.
- ...to improve your user experience and the overall quality of our services.
- When showing you tailored ads, we will not associate an identifier from cookies or similar technologies with sensitive categories.

SOME EXAMPLES OF USE OUTSIDE OF DATA HOLDERS...



MailOnline

Home News U.S. | Sport | TV&Showbiz | Australia | Femail | Health | Science | Money | Vic

Latest Headlines | News | World News | Arts | Headlines | France | Pictures | Most read | Wires | Discounts

Teacher sacked for posting picture of herself holding glass of wine and mug of beer on Facebook

By DAILY MAIL REPORTER
UPDATED: 23:45 BST, 7 February 2011





METRO

Thieves steal \$1M in jewels from Jerry Seinfeld's ex-girlfriend

By **Larry Celona** and **Joe Tacopino**

July 20, 2016 | 3:39am



Cops investigate a burglary at the home of Jerry Seinfeld's ex Shoshanna Lonstein Gruss.



Going away? Don't tell your Facebook friends or risk having your insurance claims rejected

- If you post photos of holidays you risk having insurance claims rejected if your home is burgled
- Allowing your Facebook page to reveal your location automatically can also be risky

By [RUTH LYTHE FOR MONEY MAIL](#)

PUBLISHED: 00:10 BST, 22 April 2015 | **UPDATED:** 09:32 BST, 22 April 2015



SHARE
PICTURE



Woman Loses Benefits After Posting Facebook Pics

By KI MAE HEUSSNER
Nov. 23, 2009

 Share with Facebook

 Share with Twitter



INDIRECT INFORMATION USE

Examples:

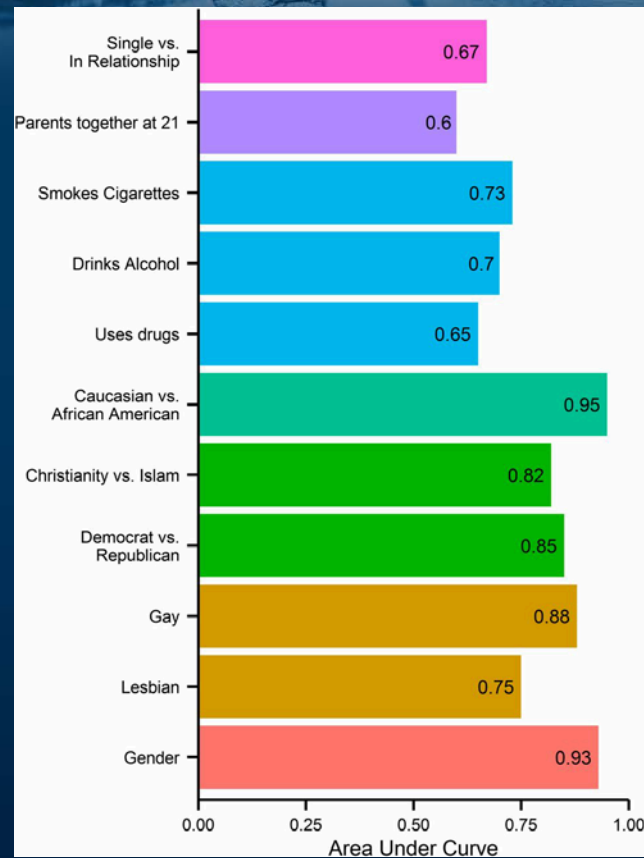
- a user who is interested in university/educational issues
 - → is very likely to be a young adult.
- in a dance club, people come together due to their common interest; in an office, people connect to each other because of similar professions
 - → to infer someone's attribute from the attributes of his/her friends.
- A martini or cigarette in your hand in 98% of your photos
 - → is very likely to get liver failure, lung cancer, and lowered life expectancy!

WHAT HAS RESEARCH SHOWN SO FAR?

Kosinski et al. (2013)

Analyzed 58,466 Facebook users :

- like history (170 likes/person on average)
- profile information
- the results of several psychometric tests



MORE EXAMPLES...

- Schwartz et al. (2013) analyzed text of 15.4 million status updates from a total of 74,941 Facebook users. Predicted gender with 92 % accuracy
- Backstrom and Kleinberg (2014) managed to predict whether a user is single or not with 68 % accuracy and whether he/she is single or married with 79 % accuracy.
- Jernigan et al. (2009) looked at sexual orientation and achieved an accuracy of 78 % by analyzing friendship associations.
- Zheleva and Getoor (2009), examined user attributes are the country, gender and political views.
- Rao et al. (2010) evaluated the accuracy of predicting gender (72 %), age (74 %), regional origin (77 %) and political affiliation (83 %) from Twitter messages.
- Conover et al. (2011) (95 % accuracy) on political views were obtained from Twitter users.
- Very good results on political views from Twitter (89 % accuracy) were also achieved by Penna et al. (2011)
- ...

MORE EXAMPLES...

- Schwartz et al. (2013) analyzed text of 15.4 million status updates from a total of 74,941 Facebook users. Predicted **gender** with 92 % accuracy
- Backstrom and Kleinberg (2014) managed to predict whether a user is **single** or **not** with 68 % accuracy and whether he/she is **single** or **married** with 79 % accuracy.
- Jernigan et al. (2009) looked at **sexual orientation** and achieved an accuracy of 78 % by analyzing friendship associations.
- Zheleva and Getoor (2009), examined user attributes are the **country**, **gender** and **political views**.
- Rao et al. (2010) evaluated the accuracy of predicting **gender** (72 %), **age** (74 %), **regional origin** (77 %) and **political affiliation** (83 %) from Twitter messages.
- Conover et al. (2011) (95 % accuracy) on **political views** were obtained from Twitter users
- Very good results on **political views** from Twitter (89 % accuracy) were also achieved by Penna et al. (2011)
- ...

THE DATABAIT TOOL

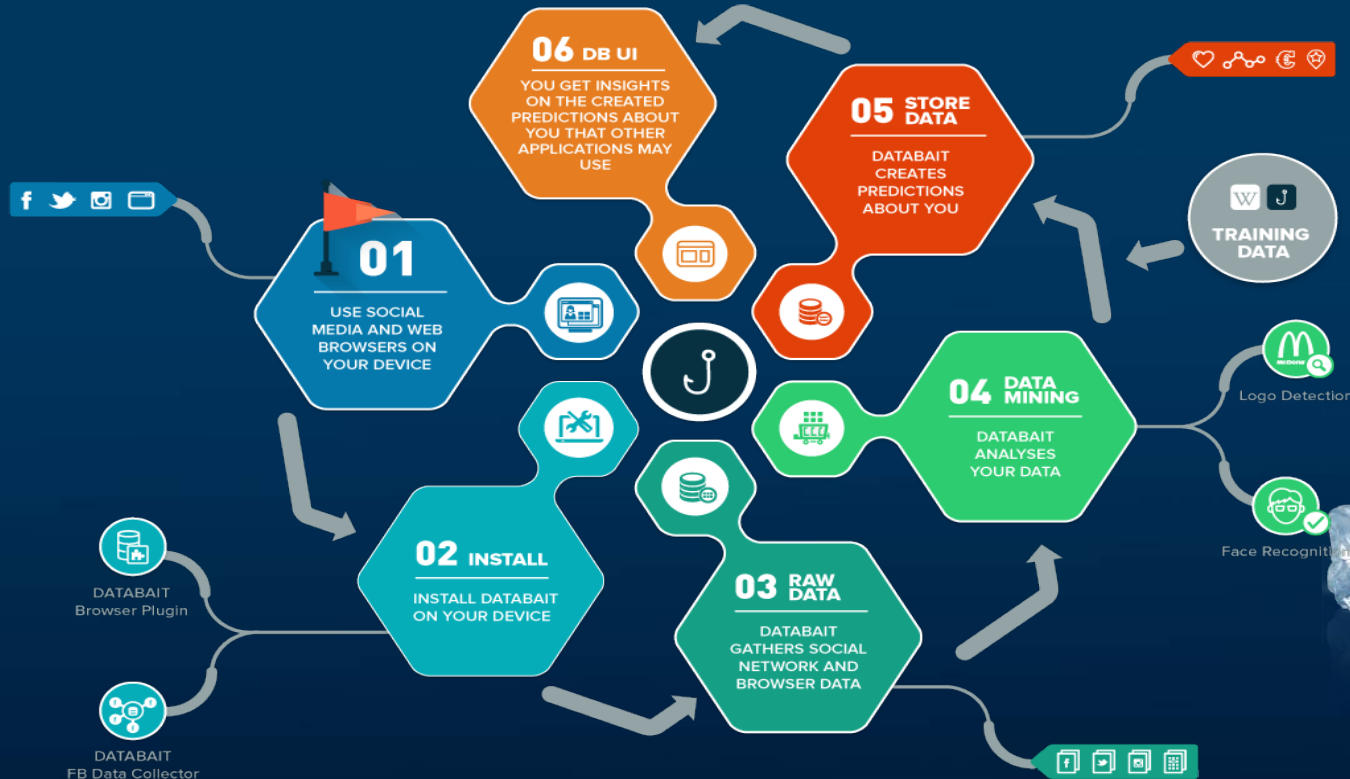
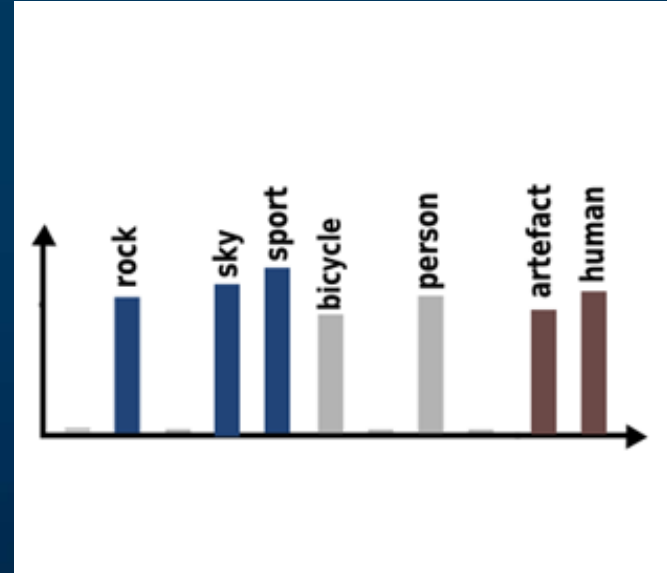


IMAGE MINING TOOL



Input



Output

IMAGE LEAKS OVER FACEBOOK



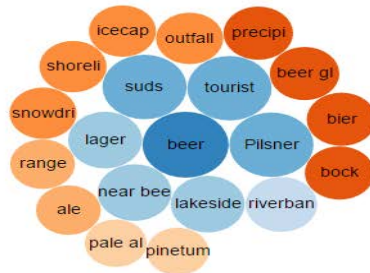
My Privacy

Overview Friends Location Leaks Image Leaks

beer Pilsner tourist suds lager near beer lakeside
riverbank bock bier beer glass precipice outfall
icecap shoreline snowdrift range ale pale ale
pinetum speaker manageress seaside ben

Retrieve my latest images

Top 20 Visual Concepts Detected



LOCATION LEAKS OVER FACEBOOK



My Privacy

[Overview](#) [Friends](#) [Location Leaks](#) [Image Leaks](#)

Cairo

Märsta Helsinki Kristiansand Oulu Stockholm

Turku

Vatican City

Weehawken

Berkeley Mikkelí

Tossa de Mar

Glasgow

Beijing

Bayswater

Tampere

Kaustinen

Rome

Varaždin

Vaduz Årsta

Bamako

Grindavík

Brasília

Chiang Mai

Kirkkonummi

Ävsbyn

Kensington

Turin

Ottawa

Rio Grande

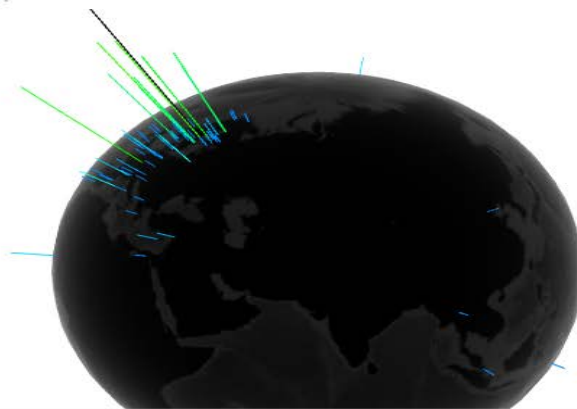
London

Kirkenes

Olderdalen

388 posts with location information

Retrieve my latest posts



EXAMPLES OF 96 TOTAL CLASSIFIED ATTRIBUTES

Health

cannabis
energydrink
coffee
nosubstance
drinkingbehavior
alcohol
smokingbehavior
cigarettes
healthstatus
bmiclass

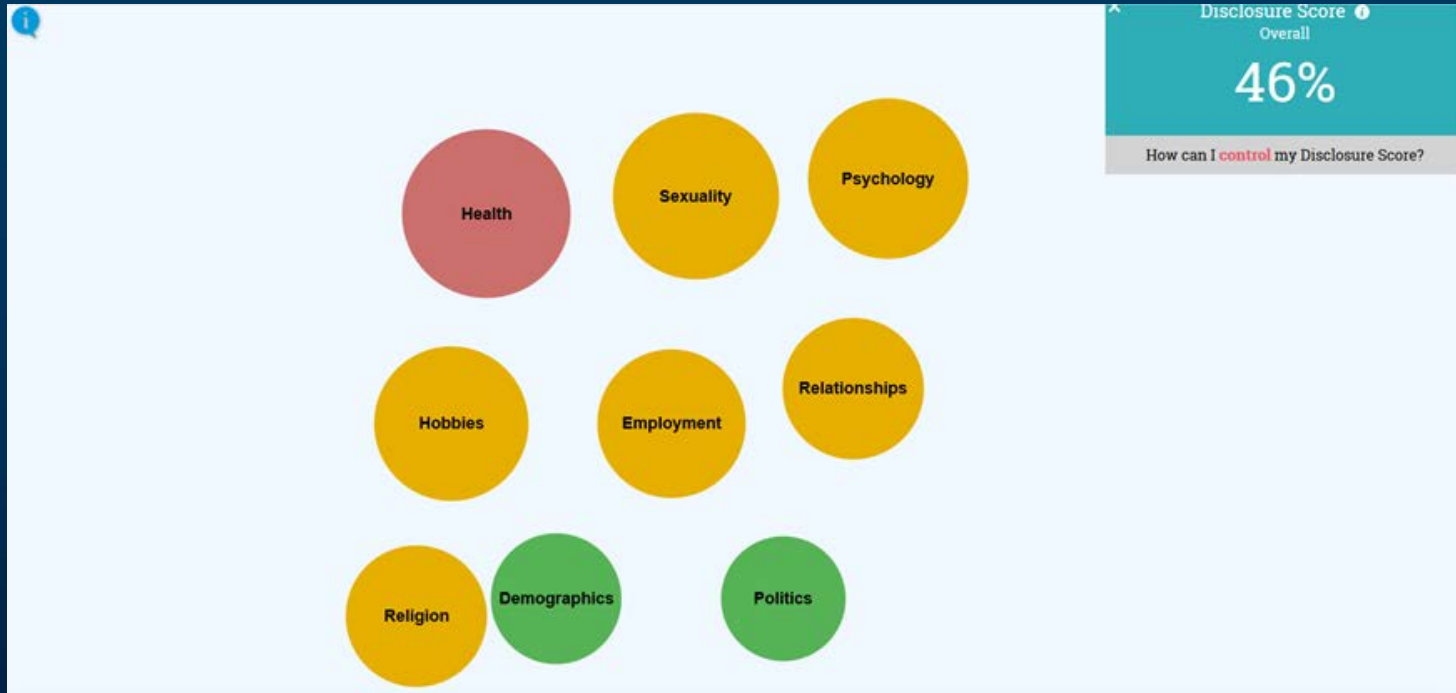
Hobbies

Swimming
Bicycling
Listening-to-music
Reading
Watching-series-or-movies-at-h
Cooking
Going-to-the-movies
Eating-out
Running

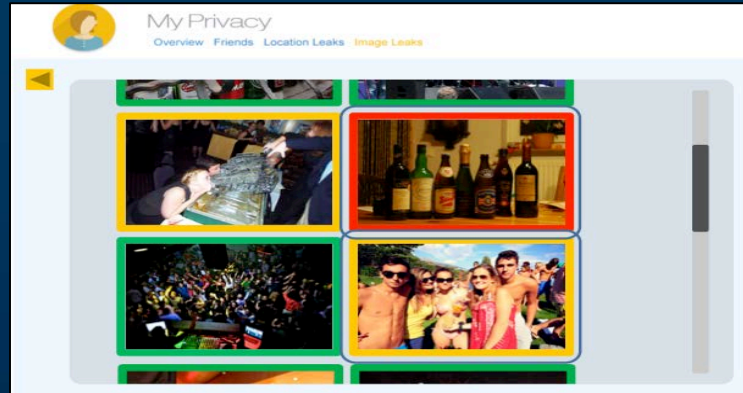
Psychology

can-be-moody
is-ingenious-a-deep-thinker
is-helpful-and-unselfish-with-
is-talkative
can-be-somewhat-careless
is-inventive
is-outgoing-sociable
is-curious-about-many-differen

PERSONAL USER ATTRIBUTES WERE ORGANIZED INTO 9 CATEGORIES



VISUALISATION PHOTO ALBUMS



My Privacy
Overview Friends Location Leaks Image Leaks

keywords
drunkard
beer drinker
cafe
brasserie beer hall

What parts of your privacy profile are affected?

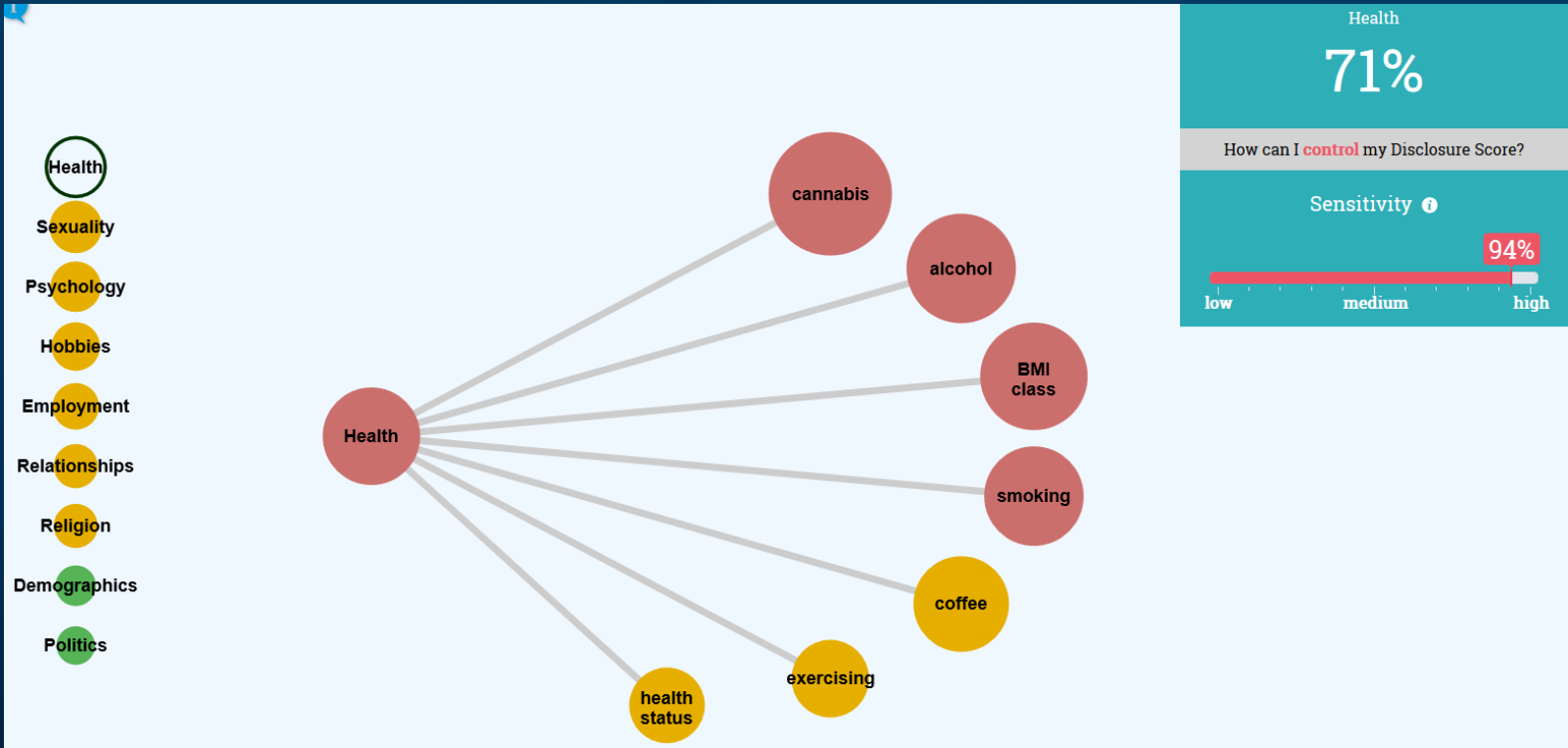
- consumer profile
- health factors

My Privacy
Overview Friends Location Leaks Image Leaks

keywords
lovers
beach wear
drunkard

What parts of your privacy profile are affected?

- sexual profile
- health factors
- consumer profile





Health

Sexuality

Psychology

Hobbies

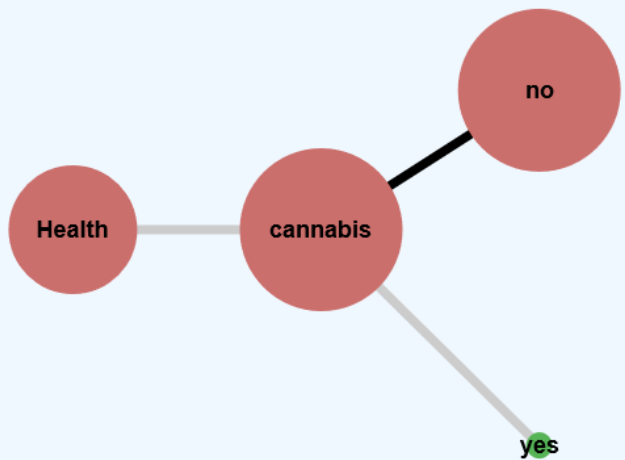
Employment

Relationships

Religion

Demographics

Politics



Disclosure Score ⓘ
Health > Cannabis

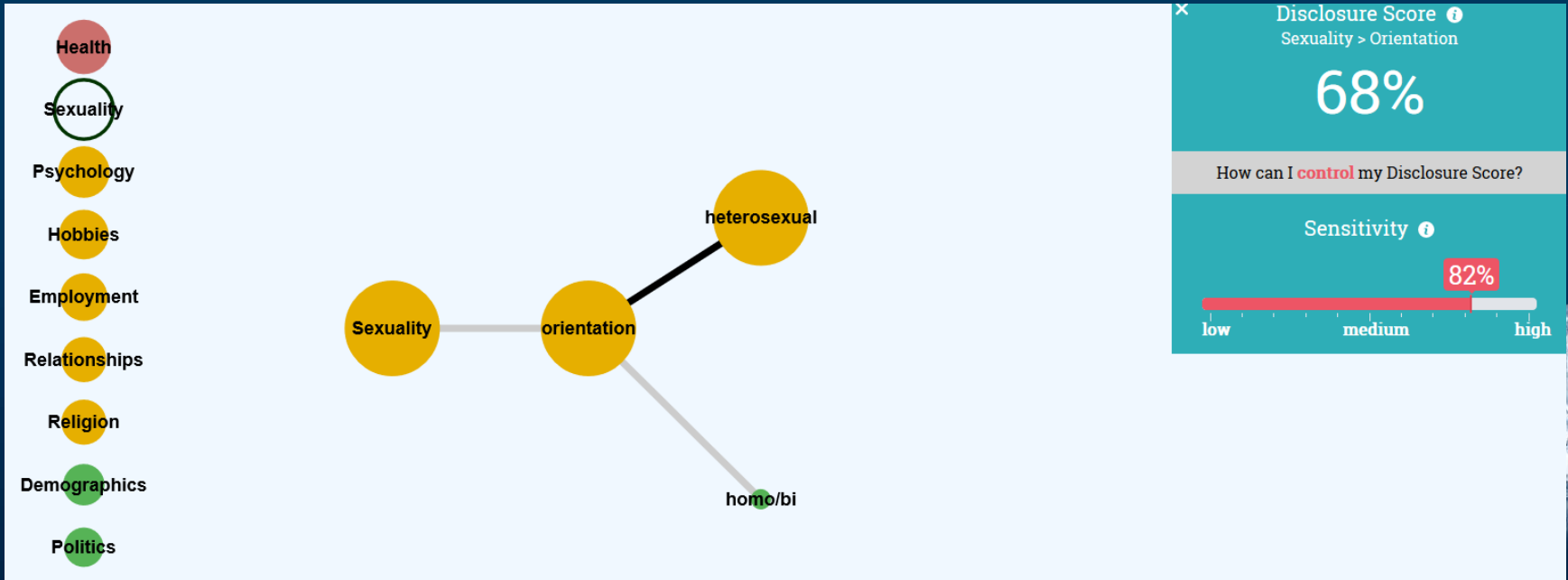
88%

How can I **control** my Disclosure Score?

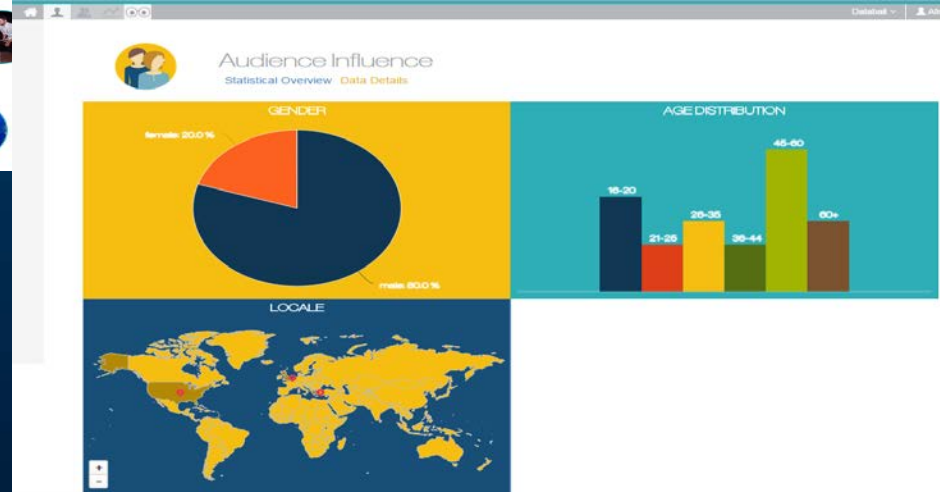
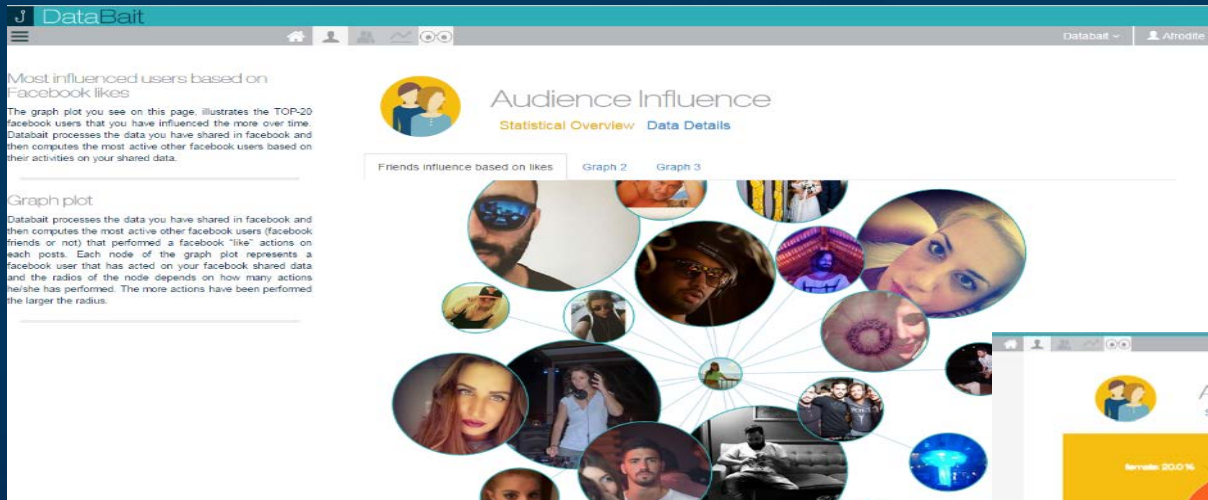
Sensitivity ⓘ

94%

low medium high



AUDIENCE INFLUENCE





OUR STUDY IN LTU

- A mixture of participants:
 - occupation (13 students, 14 non-students)
 - educational background (5 high school, 11 BA and 11 MA level)
 - gender (15 male, 12 female)
 - cultural backgrounds (14 Swedish, 13 non-Swedish)
 - age (from 18 to 58 years old)



SOME QUOTES FROM OUR PARTICIPANTS...

- *“I was recently looking for a job and I want to be very careful about what others tag me in and what I post about myself ... if I knew that, then I would go and make changes so that I get things that only I think are ok to see.”*
- *“[Facebook] has all kinds of crazy programs and algorithms that analyze all your habits and likes. Facebook probably does worse than the developers of DataBait on a regular basis. Plus they are a massive multi-national company but DataBait is just a couple of developers in Europe.”*
- *“When you share your photos you share maybe, one, two, or three photos, and you kind of forget the ones you previously posted, and then when you see them all together, it gives you a kind of summary of the pictures that you are introducing of yourself, the profile that you are actually producing.”*



IMPLICATIONS FOR DIGITALIZATION

- If in this project we can infer, then others can do it too!
- How much are we aware of data gathered related to end users?
- How much personal information is susceptible to data mining?
- Indirect information disclosure
- User cognitive ability lags behind technological advances
- How can we empower users regarding data mining on the information they are sharing?

THANK YOU!

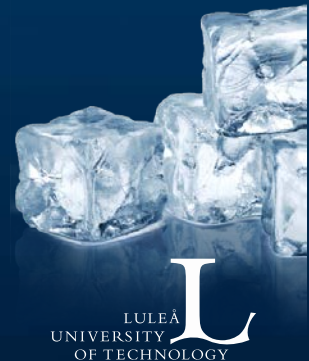


Ali.Padyab@ltu.se



[linkedin.com/in/alipadyab](https://www.linkedin.com/in/alipadyab)

Acknowledgments





LULEÅ
UNIVERSITY
OF TECHNOLOGY

