

CORSO DI
WEB MINING E RETRIEVAL
- INTRODUZIONE AL DEEP LEARNING -

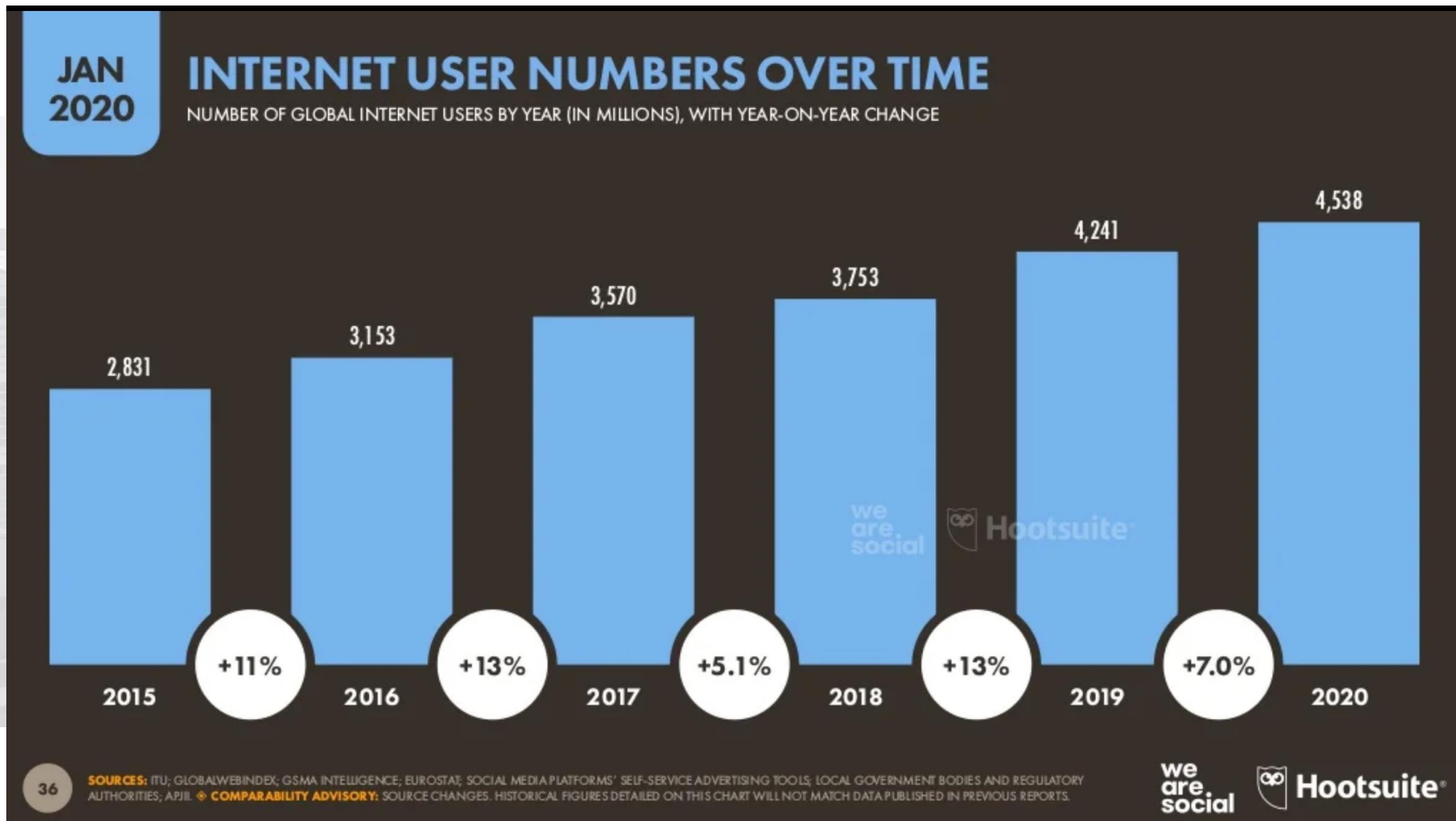
Corso di Laurea in Informatica, Ing. Gestionale, Ing.
Internet, Ing. Informatica,
(a.a. 2022-2023)

Roberto Basili

Overview

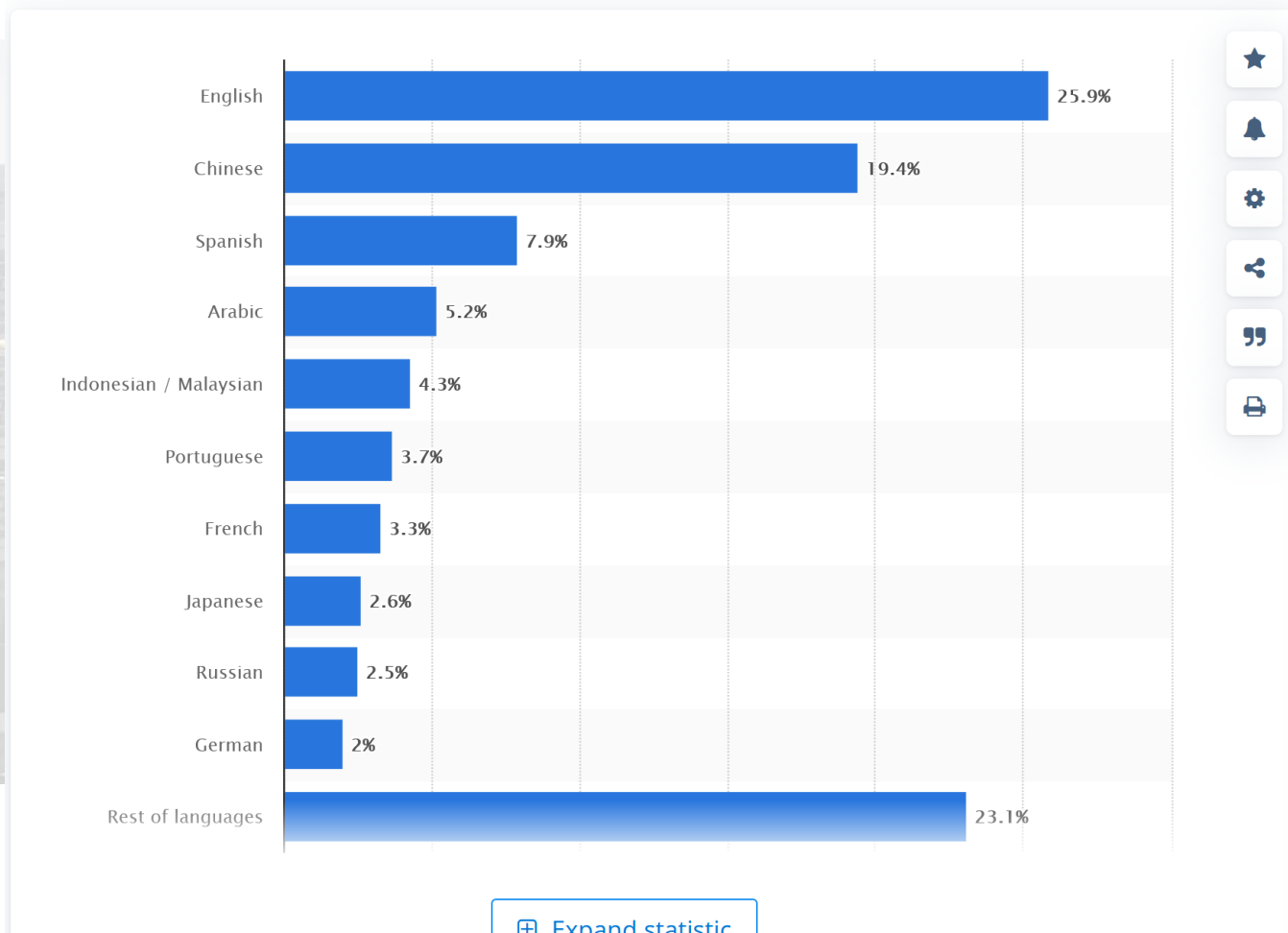
- Web Mining & Retrieval: Motivations & perspectives
 - Web, User-generated contents, Social Media
 - The role of *learning*
 - What is Machine Learning?
 - Data-driven algorithms: sources of complexity
- Main Applications
 - Intelligent Web Search
 - User Profiling for Marketing or Brand reputation management
 - Web Recommending
 - Spoken Dialogue Interaction in Robotics or in Web/mobile Interfaces

Internet statistics (Jan 2021)

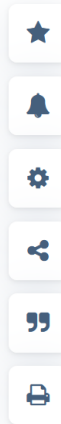
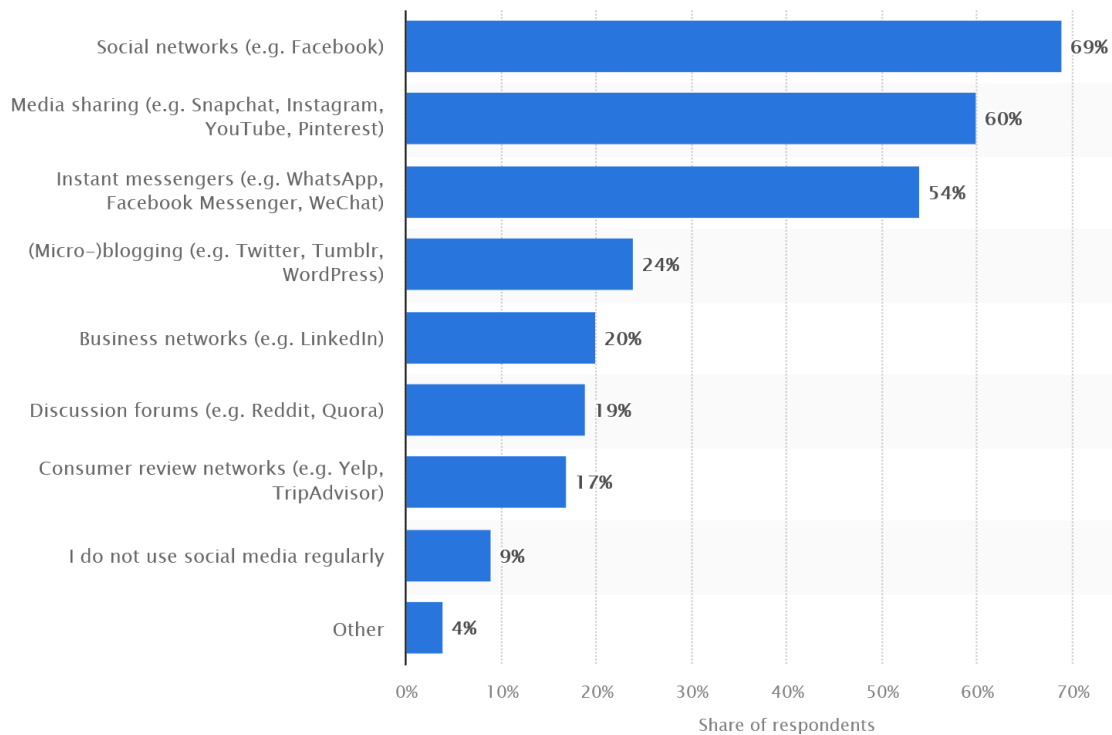


Internet Statistics (Jan 2020)

Most common languages used on the internet as of January 2020, by share of internet users



What kinds of social media do you use regularly?



© Statista 2022 🇩🇪

Details:

[Show source](#) ⓘ

Do you know

More than
4,000 new books
are published every day



Do you know

**Contains more
information than a
person was likely to
come across
in a lifetime in the
18th century...**



**JAN
2022**

ESSENTIAL DIGITAL HEADLINES

OVERVIEW OF THE ADOPTION AND USE OF CONNECTED DEVICES AND SERVICES



GLOBAL OVERVIEW

TOTAL
POPULATION



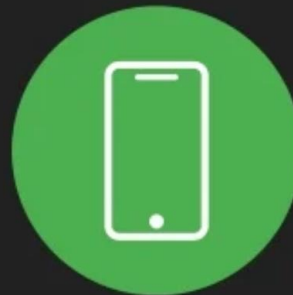
we
are
social

**7.91
BILLION**

URBANISATION

57.0%

UNIQUE MOBILE
PHONE USERS



**5.31
BILLION**

vs. POPULATION

67.1%

INTERNET
USERS



**4.95
BILLION**

vs. POPULATION

62.5%

ACTIVE SOCIAL
MEDIA USERS



**4.62
BILLION**

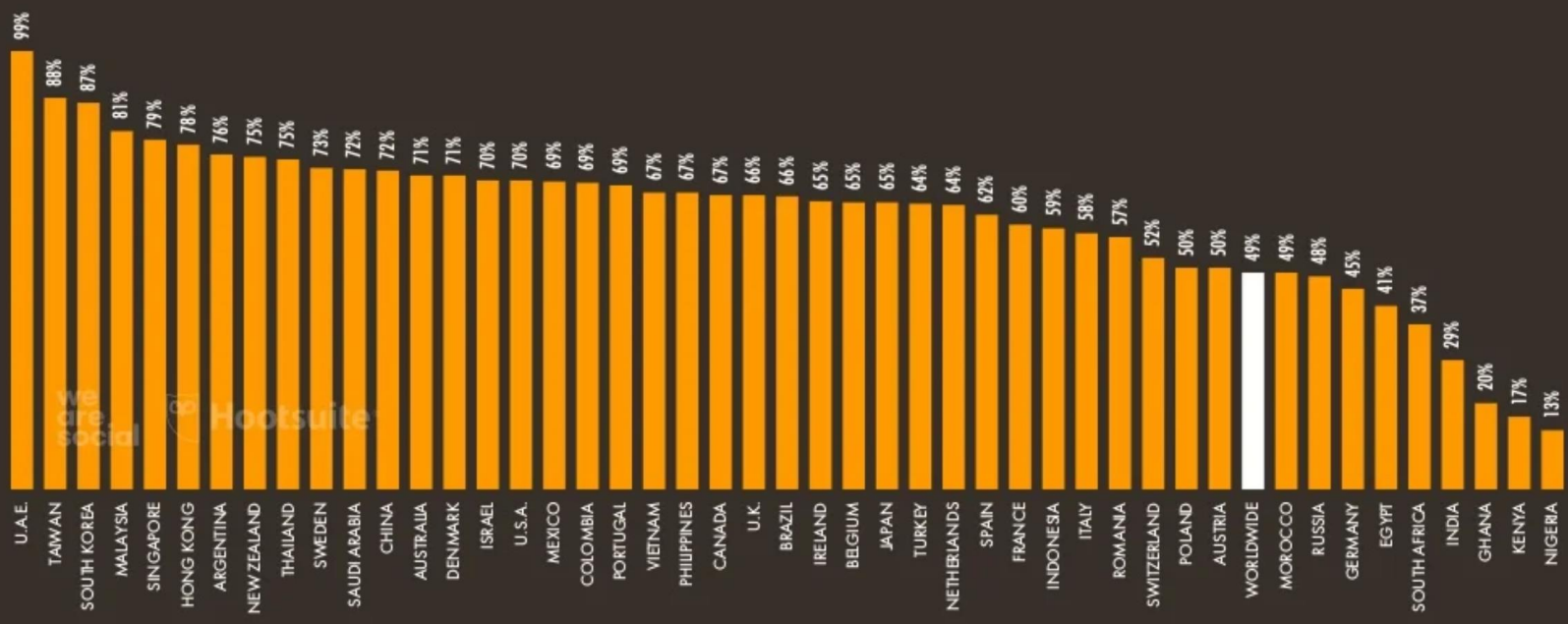
vs. POPULATION

58.4%

JAN 2020

SOCIAL MEDIA PENETRATION

THE NUMBER OF ACTIVE SOCIAL MEDIA USERS COMPARED TO TOTAL POPULATION, REGARDLESS OF AGE



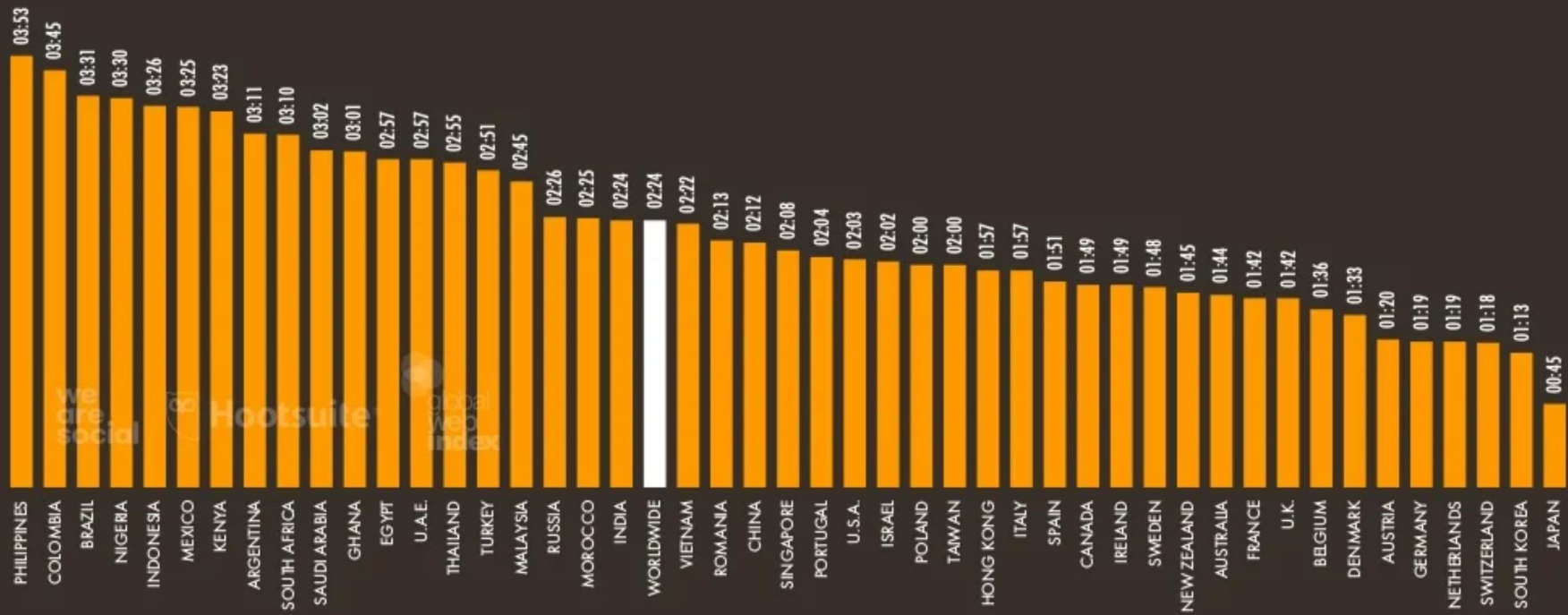
81

SOURCES: KEPIO ANALYSIS; COMPANY STATEMENTS AND EARNINGS ANNOUNCEMENTS; SOCIAL MEDIA PLATFORMS' SELF-SERVICE ADVERTISING TOOLS; MEDIASCOPE; CAFEBAZAR (ALL LATEST DATA AVAILABLE IN JANUARY 2020). ***NOTES:** PENETRATION FIGURES ARE FOR TOTAL POPULATION, REGARDLESS OF AGE. **◆ COMPARABILITY ADVISORY:** SOURCE AND BASE CHANGES.

JAN
2020

DAILY TIME SPENT USING SOCIAL MEDIA

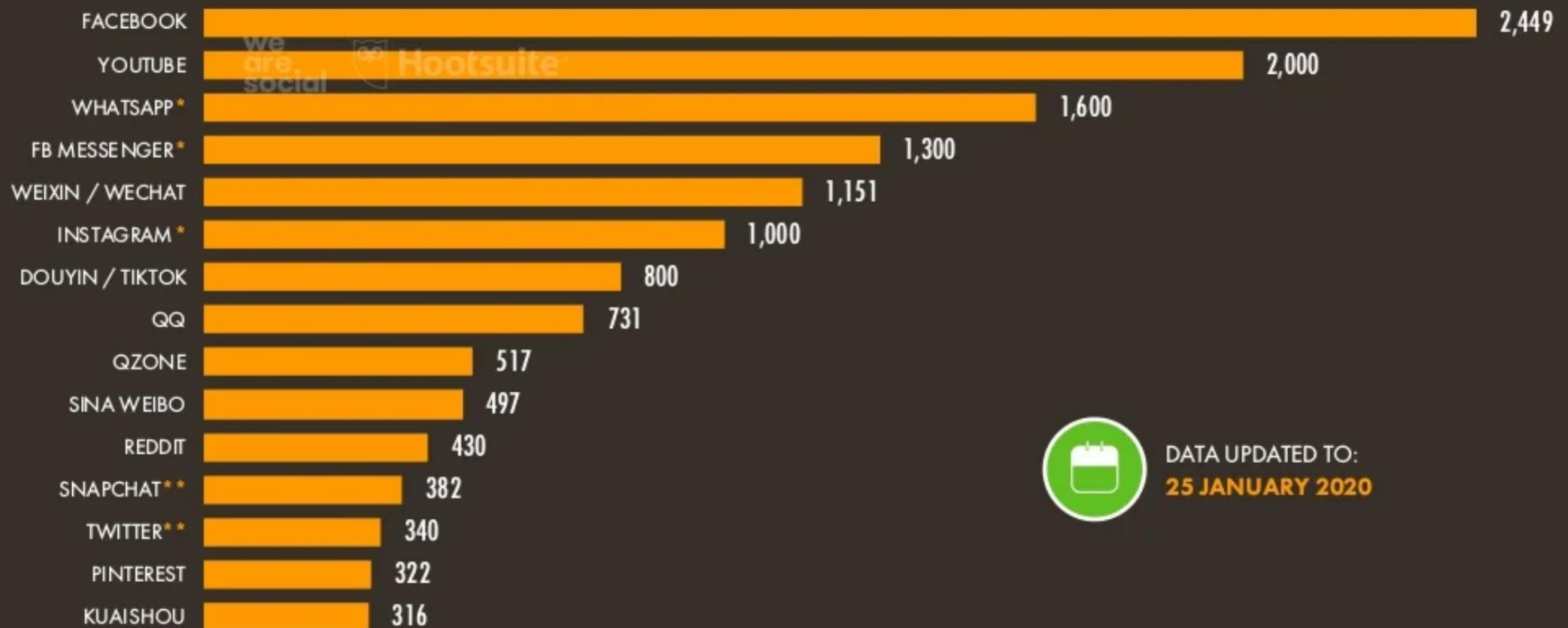
AVERAGE DAILY TIME (IN HOURS AND MINUTES) THAT INTERNET USERS AGED 16 TO 64 SPEND USING SOCIAL MEDIA ON ANY DEVICE



**JAN
2020**

THE WORLD'S MOST-USED SOCIAL PLATFORMS

BASED ON MONTHLY ACTIVE USERS, ACTIVE USER ACCOUNTS, ADVERTISING AUDIENCES, OR UNIQUE MONTHLY VISITORS (IN MILLIONS)



DATA UPDATED TO:
25 JANUARY 2020

**JAN
2022**

MAIN REASONS

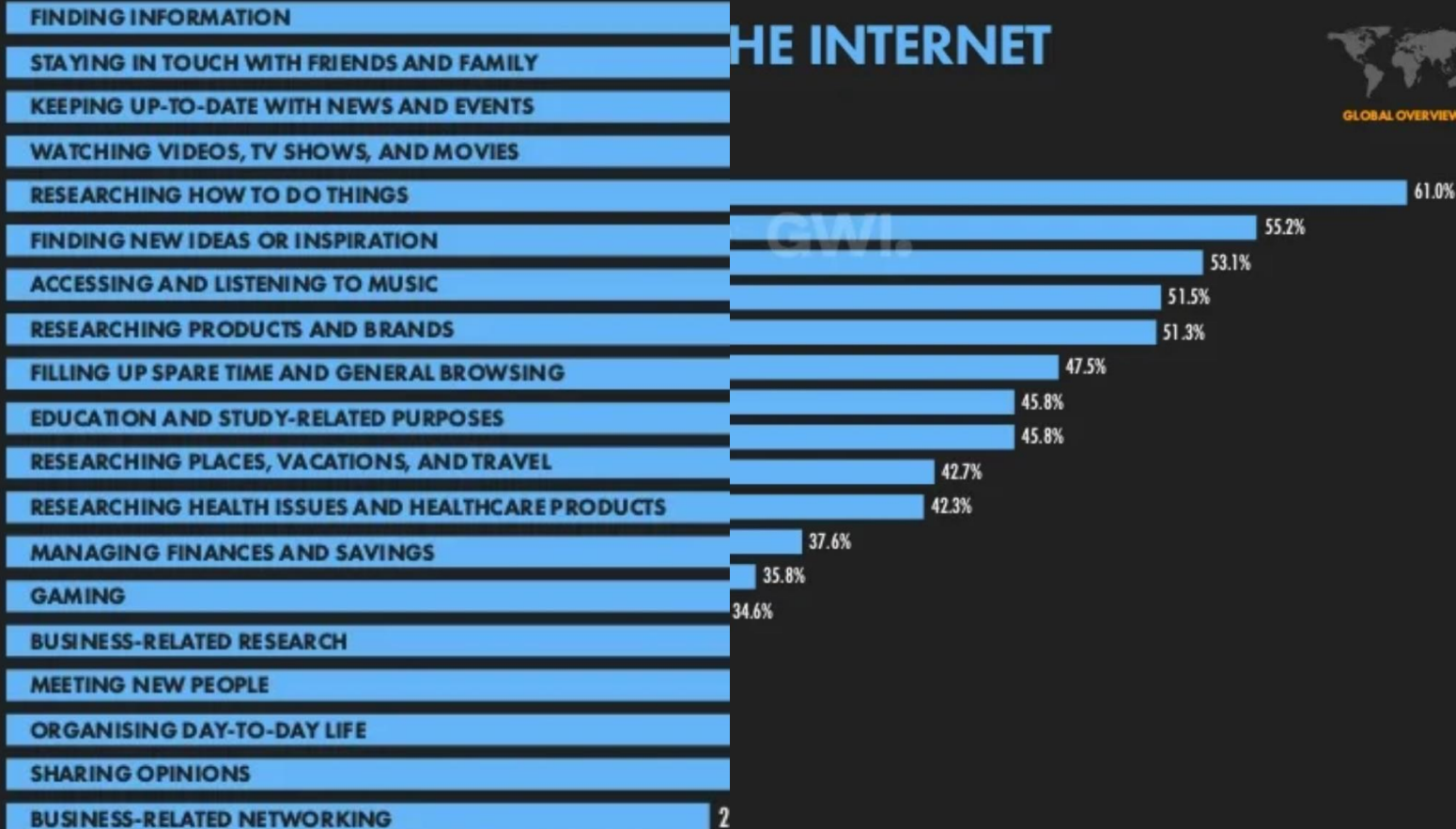
PRIMARY REASONS WHY INTERNET USERS AGED

**JAN
2022**

THE INTERNET



FINDING I
STAYING I
KEEPING U
WATCHIN
RESEARCH
FINDING I
ACCESSIN
RESEARCH
FILLING U
EDUCATIO
RESEARCH
RESEARCH
MANAGIN
GAMING
BUSINESS
MEETING I
ORGANIS
SHARING
BUSINESS



JAN
2020

TWITTER AUDIENCE OVERVIEW

THE POTENTIAL NUMBER OF PEOPLE THAT MARKETERS CAN REACH USING ADVERTS ON TWITTER

NUMBER OF PEOPLE THAT
TWITTER REPORTS
CAN BE REACHED WITH
ADVERTS ON TWITTER



339.6
MILLION

SHARE OF POPULATION
AGED 13+ THAT MARKETERS
CAN REACH WITH
ADVERTS ON TWITTER



5.6%

QUARTER-ON-
QUARTER CHANGE
IN TWITTER'S
ADVERTISING REACH



-3.1%

PERCENTAGE OF
ITS AD AUDIENCE
THAT TWITTER
REPORTS IS FEMALE*



38%

PERCENTAGE OF
ITS AD AUDIENCE
THAT TWITTER
REPORTS IS MALE*



62%

**JAN
2020**

MOST-USED EMOJI ON TWITTER

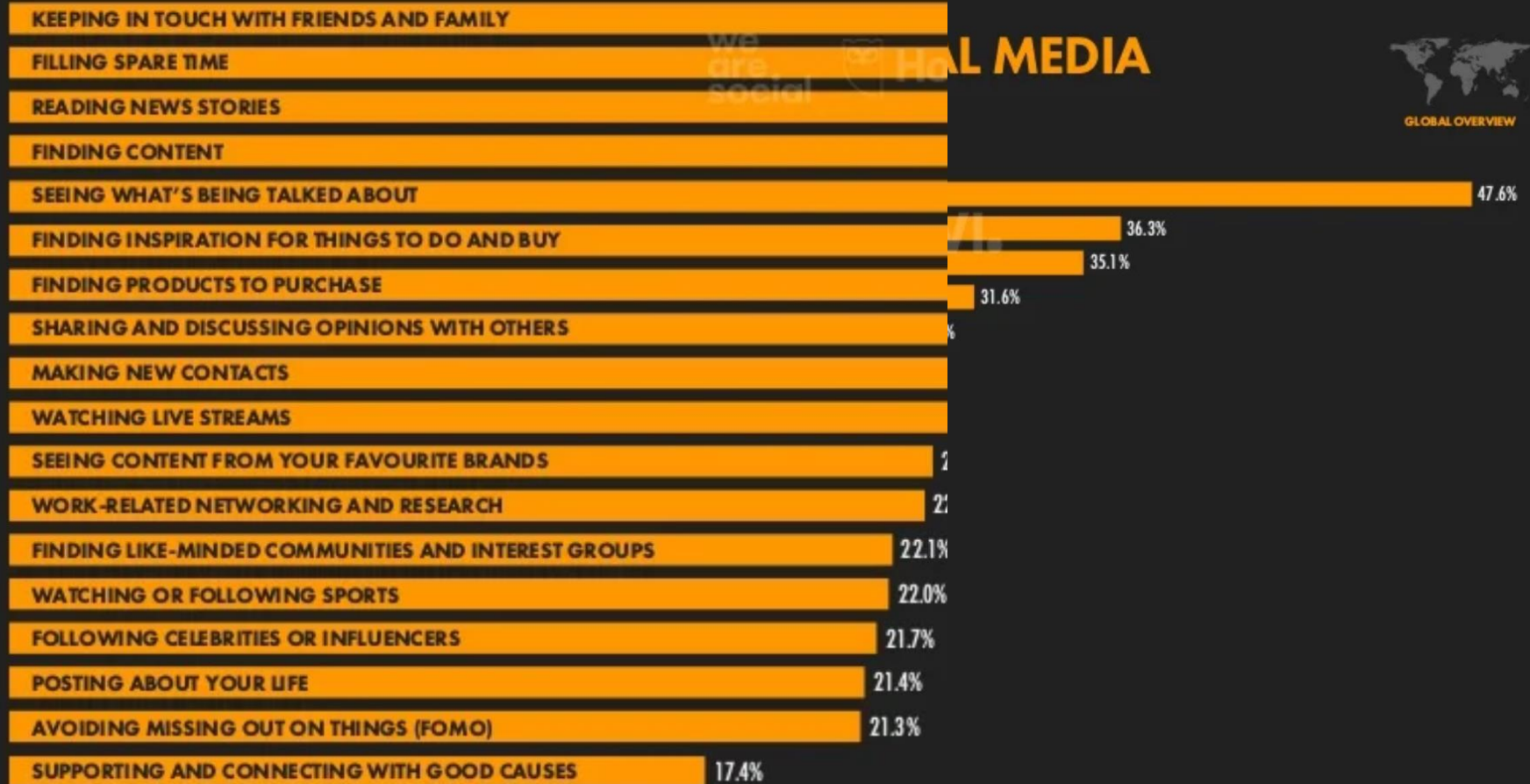
EMOJI THAT HAVE BEEN USED THE GREATEST NUMBER OF TIMES ON TWITTER (ALL TIME)

#	EMOJI	TIMES USED	#	EMOJI	TIMES USED	#	EMOJI	TIMES USED	#	EMOJI	TIMES USED
01		2,671,000,000	11		428,000,000	21		245,000,000	31		198,000,000
02		1,289,000,000	12		389,000,000	22		238,000,000	32		193,000,000
03		966,000,000	13		382,000,000	23		237,000,000	33		191,000,000
04		964,000,000	14		365,000,000	24		236,000,000	34		187,000,000
05		817,000,000	15		359,000,000	25		232,000,000	35		182,000,000
06		743,000,000	16		336,000,000	26		229,000,000	36		181,000,000
07		632,000,000	17		309,000,000	27		217,000,000	37		168,000,000
08		500,000,000	18		273,000,000	28		216,000,000	38		165,000,000
09		493,000,000	19		258,000,000	29		212,000,000	39		163,000,000
10		475,000,000	20		246,000,000	30		199,000,000	40		163,000,000

**JAN
2022**

MAIN REASONS FOR USING SOCIAL MEDIA

PRIMARY REASONS WHY INTERNET USERS AGED 16 TO 64 USE SOCIAL MEDIA



we
are
social

GLOBAL MEDIA



GLOBAL OVERVIEW

WE ARE SOCIAL'S PERSPECTIVE: SOCIAL IN 2020

SHIFTS IN HOW PEOPLE BEHAVE AND INTERACT ON SOCIAL



BAD INFLUENCE

Being a creator has lost its lo-fi sheen; many lifestyle influencers lead unrelatable lives, while celebrity 'creators' like Will Smith are blowing up on platforms like YouTube and TikTok. As a result, there's a growing backlash against influencer culture and the metrics that drive it.

In 2020, brands will look beyond likes, followers and reach to generate genuine engagement



ADDED VALUE

The internet has long been a wild west where intellectual property is barely there. But in a maturing digital frontier, creators have grown dedicated audiences who not only see value in their content, but recognise their style anywhere. As a result, communities are rallying to protect creators.

In 2020, brands will take greater steps to ensure they're being respectful of digital communities



RUNNING COMMENTARY

Audiences are increasingly willing to invest time and attention in content and narratives they deem to have a higher value. This isn't about a shift back to traditional media. It's about longer, more complex content designed to be consumed in-platform and on smaller screens.

In 2020, brands will tell more complex stories across multiple touchpoints on social

Dealing with *real* Social media data



WM&R: Motivations

- *What does Web Mining mean?*
- *Why Information Retrieval is involved?*
- *Why Machine Learning and mostly Deep Learning?*
- *Which are the contributions of IR/ML/NLP to technologies that support and exploit Web Data, Information and Knowledge?*
- *Which are the technological perspectives in the medium-long term?*

What is Web Mining?

- *Web Mining* refers to a body of technologies currently needed for the *exploitation of publicly available information from the Web and the IoT*
 - Contents: data but also ... PEOPLE, LOCATIONS, EVENTS, CONCEPTS, TEMPORAL INFORMATION ...
 - Relations:
 - Links within structured networks (retweets, follows, ...)
 - Thematic, interpersonal and semantic associations
 - Similarities and Analogies among people, behaviours, preferences
 - On-Line Structured and semi-structured resources (e.g. Wikipedia)
 - Textual, Multimedia and Multilingual Contents
 - Trends e time-related information (community on-line behaviours)
 - Opinions, Preferences, Expectations

Why IR?

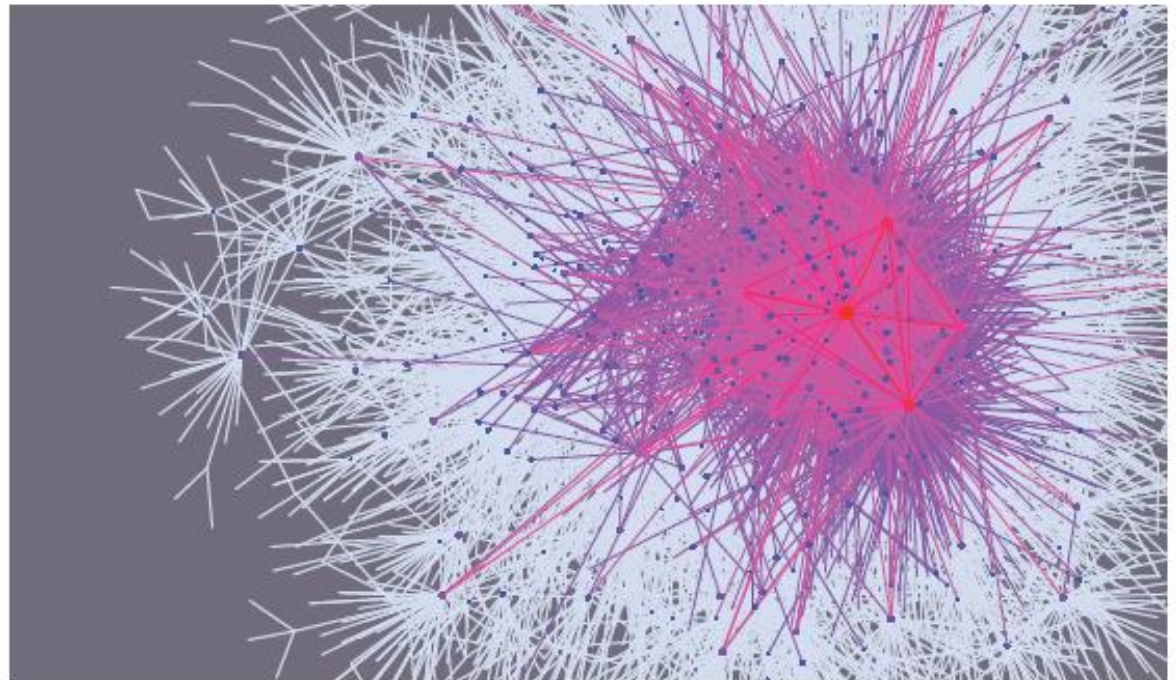
- The volumes involved in Web Mining pose the crucial problem of *locating information* beforehand
- Automatic information access is possible only if we solve the two major challenges
 - **What** is relevant
 - **Where** the relevant information is located
- **Searching information corresponds to computing an uncertain function that models the mapping between information needs and the targeted data**

Machine Learning vs IR?

- Web mining involve heterogeneous information that is characterized search as strongly uncertainty process
- The available information is characterized by:
 - Incompleteness:
 - Short queries as an incomplete description of the information need
 - Variability: Wealth of data vs. heterogeneity of formats and access modes
 - Contents are dispersed in various forms across data sources
 - Vague Requirements
 - Information is often implicit (i.e. partially and qualitatively expressed) in the operational contexts
 - Subjectivity
 - Relevance depends on the user and not just on the contents
 - Timeliness
 - Authority

Machine Learning vs. IR

- Uncertainty is so pervasive that exhaustive solutions (i.e. global *optima*) are not available or even not existing
- “*Finding diamonds in the rough*” (Fan Chung, UCSD)



Machine Learning vs. IR

- ML technologies offer a wide variety of algorithms, strategies and techniques for the induction of sub-optimal, but surprisingly effective, solutions from available data
- Through *learning* data can be effectively used to suggest retrieval hypothesis, that are models of the *mapping* function (Learning to search)
- What is the target of the learned function? To improve computational aspects of the currently applied processes, such as
 - Semantic Accuracy (i.e. best answers first)
 - System Responsiveness (i.e. reducing speed of the retrieval process)
 - Resource usage (i.e. more effective with less memory or input data)

Machine Learning

- Machine learning is **the study of computer algorithms that allow computer programs to automatically improve through experience**. (Tom Mitchell, *Machine Learning*, McGraw-Hill, 1997)
- The evidence of the success of a learning process corresponds to the possibility of observing a measurable increment ΔP of performances in solving a task C on the basis of experiences E that the agent is able to gather during its lifecycle.
- The nature and complexity of the learning ability is fully confined to the ability of characterizing the primitive notions here involved:
 - TASK C
 - PERFORMANCE P
 - EXPERIENCE E

Experience and Learning

- Forms of experiences
- In chess games:
 - Data on previous matches, such as won challenges (or defeats) able to gather the utility (or inadequacy) of the strategies or moves there carried out.
 - Evaluation about individual moves offered by an external teacher (oracle, guide).
 - Adequacy of individual behaviour derived from self-observation, such as the capability of analysing matches against itself based on an existing explicit model of the rules and strategies of the game.

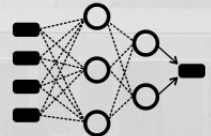
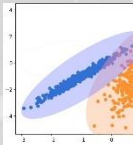
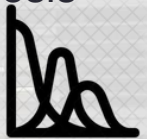
ML: a visual introduction

- See URL: http://www.r2d3.us/visual-intro-to-machine-learning-part-1/?imm_mid=0d76b4&cmp=em-data-na-na-newsltr_20150826

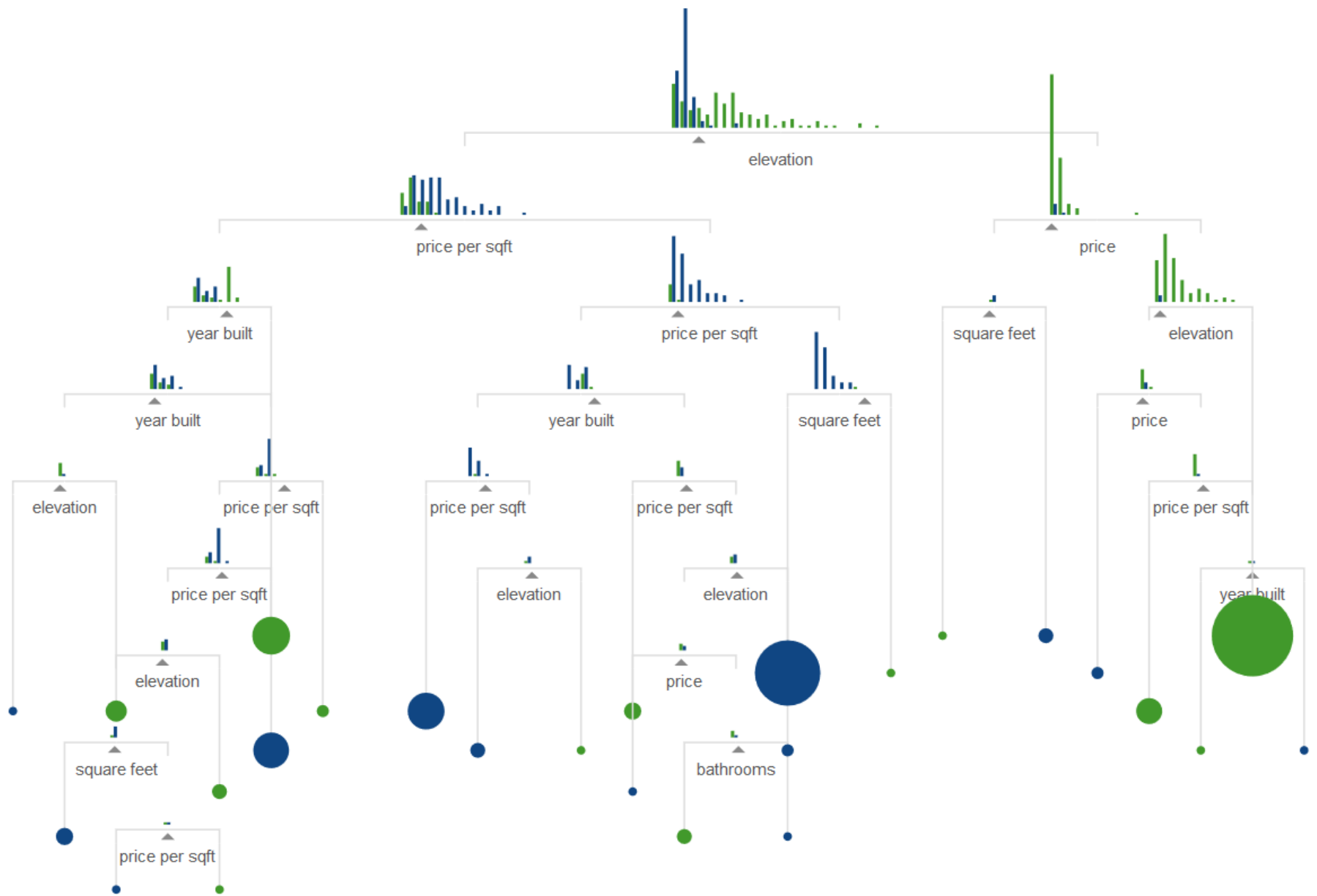


The mathematics of Learning

- Learning corresponds to the induction of mathematical function (i.e. the decision rules) that may have a discrete as well as a continuous behaviour:
 - Logical functions, (ad es., decision trees)
 - Learning the rules that better explain the data
 - Induction: Recursive search for necessary and sufficient conditions.
 - Probabilistic Approaches:
 - Learning what is *most likely* to be the better decision, according to an hypothesis about the input distribution (e.g., Bayesian classification)
 - Induction: Estimate the Posterior Probability (as parameters of known laws).
 - Metric Approaches
 - Decision as discrimination in metric spaces (e.g. linear and non linear functions)
 - K-NN
 - Linear Classifiers, perceptrons, Neural Networks, Support Vector Machines,...
 - Modeling as vectorial embedding, spectral analysis (space transformations)
 - Induction: determine the optimal parametrization from specific function classes (e.g. multilayer networks, polynomials of degree n)



Es. Decision Tree Learning



Unsupervised Learning

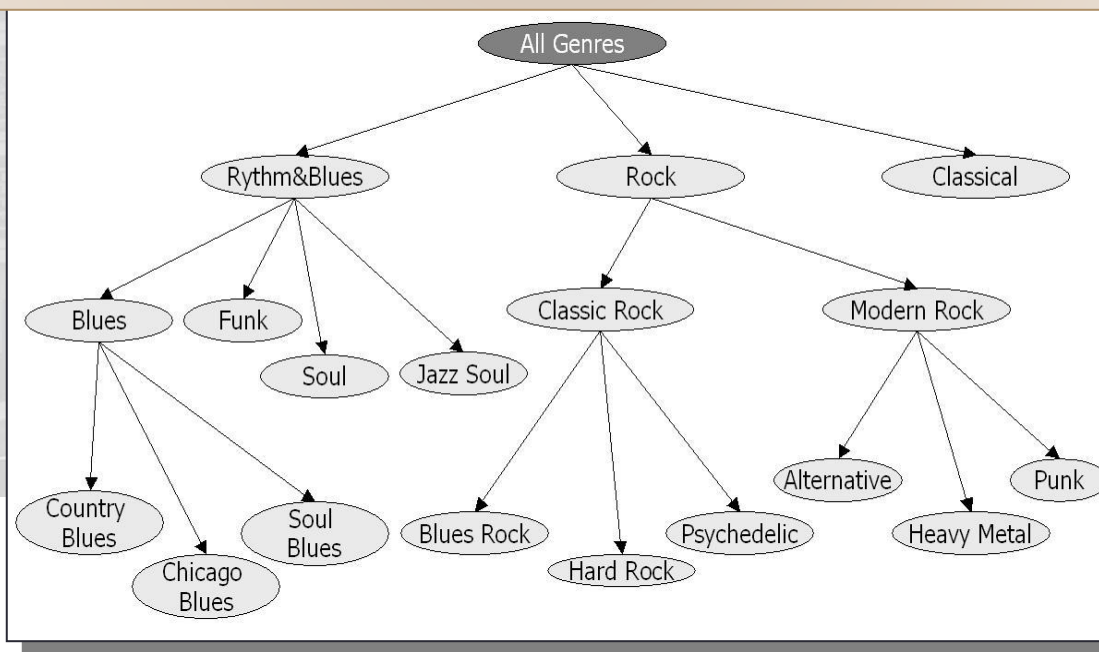
- When no oracle or knowledge of the task is available, learning may still be applied in several approaches:
 - Improve the current world model (*knowledge acquisition/discovery*)
 - Improve the efficiency of the currently available algorithms, through computational optimization
 - Better data structures representing the problem and the domain
 - Reduction of the processing steps required by the current models

Unsupervised Learning

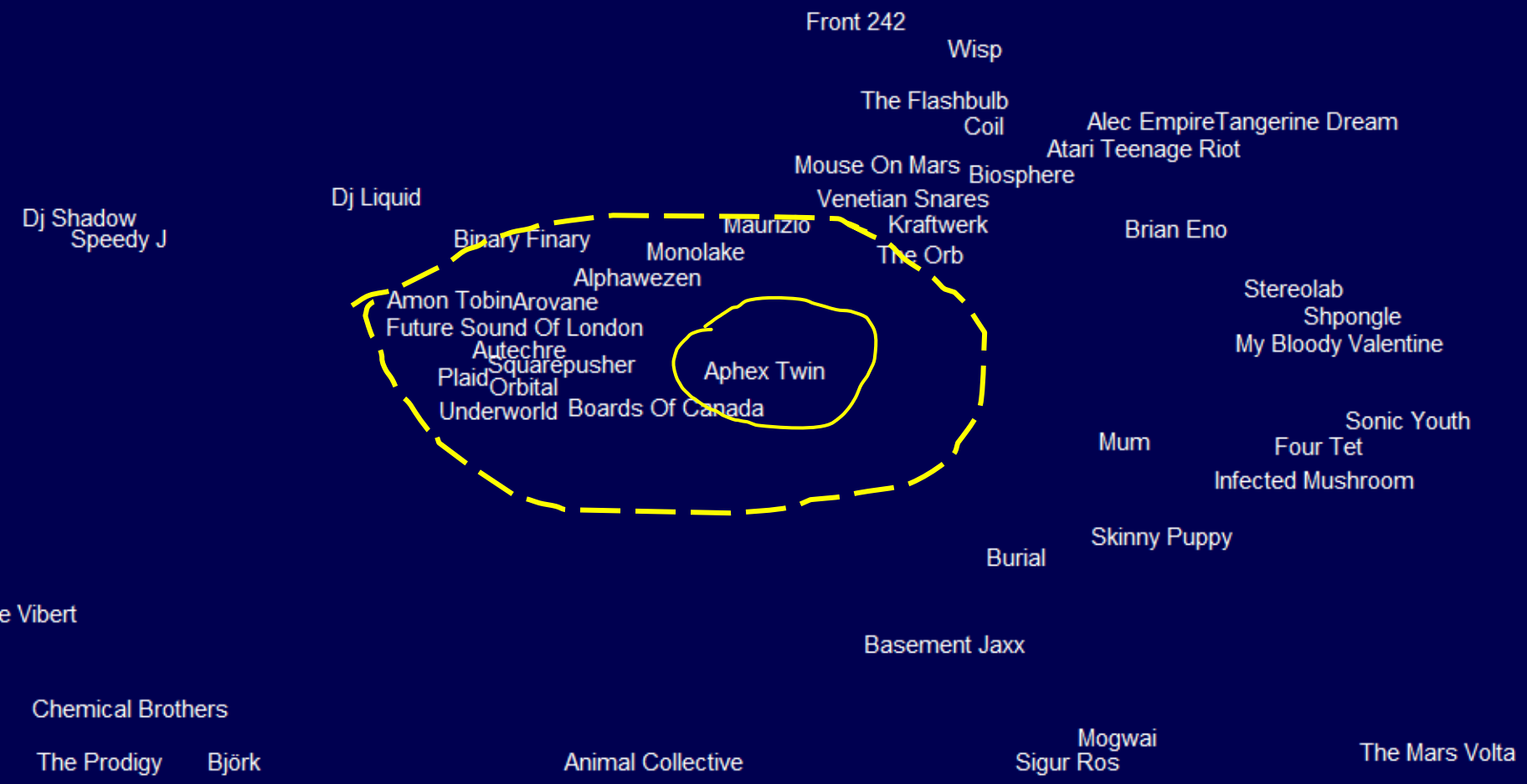
- The induced hierarchical model expresses a system of classes and relations able to improve future interaction with the song collection

It has been discovered from data

No top-down design has been applied, as in knowledge engineering, but only bottom-up inferences (i.e. generalization from data)



map



Information, Web and language

Hu meets KMT honorary chairman in Hawaii - People's Daily Online - Mozilla Firefox

File Modifica Visualizza Cronologia Segnalibri Yahoo! Strumenti Aiuto

Hu meets KMT honorary chairman

Indietro Avanti Download

Chinese President Hu Jintao (R) shakes hands with Honorary Chairman of the Chinese Kuomintang (KMT) Lien Chan, in Honolulu, Hawaii, the U.S., Nov. 11, 2011. (Xinhua/Huang Jingwen)

HONOLULU, United States, Nov. 11 (Xinhua) -- Hu Jintao, general secretary of the Central

Latest News: • Indonesia to host European Higher Education Fair

Beijing Sunny 15 / 1 City Forecast

Home >> China Politics

Hu meets KMT honorary chairman in Hawaii

(Xinhua)

11:10, November 12, 2011 🔍 +-



Chinese President Hu Jintao (R) shakes hands with Honorary Chairman of the Chinese Kuomintang (KMT) Lien Chan, in Honolulu, Hawaii, the U.S., Nov. 11, 2011.

Selections for you



Miao ethnic group celebrates Miao's New Year in SW China



World's first Angry Birds exclusive shop opens in Helsinki

Who is Hu Jintao?

Most Popular

- 1 Hu reaffirms support to Hong Kong's sta...
- 2 Hu meets KMT honorary chairman in Hawaii
- 3 China in APEC: a mutually beneficial en...
- 4 Night life in Shanghai
- 5 China's 2011 foreign trade to grow 20 p...
- 6 Beijing house prices stumble 5.1 pct as...
- 7 Lama students start school in Tibet Col...
- 8 Police in central China crack phoney ca...



Hu Jintao



Ricerca

Circa 725.000 risultati (0,09 secondi)

- Tutto
- Immagini
- Mappe
- Video
- Notizie
- Shopping
- PIÙ cont...

Tutti i ri
Per argomento

- Qualsiasi dimensione
- Grandi
 - Medie
 - Icone
 - Maggiori di...
 - Dimensioni esatte...

- Qualsiasi colore
- A colori
 - Bianco e nero
-

- Qualsiasi tipo
- Volti
 - Foto
 - Clip art
 - Disegni

Visual standard
Mostra dimensioni



Content Semantics and Natural Language

- Human languages are the main carrier of the information involved in processes such as *retrieval*, *publication* and *exchange* of knowledge as it is associated to the open Web contents
- Words and NL syntactic structures express concepts, activities, events, abstractions and conceptual relations we usually share through data
- “*Language is parasitic to knowledge representation languages but the viceversa is not true*” (Wilks, 2001)

Semantics and News

Applicazioni Risorse Sistema mar 27 lug, 23.47 dan

Gmail ... x SRL_EN x Come ... x R Econo... x Googl... x Tanl It... x Frame... x SRL_EN x Econo... x

file:///home/danilo/Downloads/SRL_ITA/sorgente/Economia%20-%20Repubblica.it.html

Telefilm in stream... Telefilm in stream... Flash Forward pri... Telefilm in stream... Cronologia Altri Pr



L'ad punta a nuove regole sulla base del modello Pomigliano. L'annuncio, che prevede l'uscita da Federmeccanica, domani al vertice con il governo o giovedì con una lettera a Bombassei. Potrebbe avvenire assieme alla decisione di creare una new company per

Pomigliano di SALVATORE TROPEA

Cisl-Uil: "L'accordo di categoria non si tocca" di S. PAROLA

Sacconi: "Su Fiat partita aperta"

Nasce Fabbrica Italia Pomigliano

Si dimette il capo di Bp buonuscita un milione di sterline



Oggi l'annuncio: a Tony Hayward subentrerà il direttore esecutivo Robert Dudley. **I costi legati al disastro sono saliti a 32,2 miliardi di dollari**, ma la società li detraerà evitando di versare al fisco Usa 10 miliardi

Manager Usa, è Ellison di Oracle il più pagato del decennio



Ha guadagnato 1,84 miliardi di dollari. Nella classifica del *Wall Street Journal* sui leader delle società quotate, secondo con 1,14 miliardi il capo di Expedia, terzo Irani di Occidental Petroleum. Solo quarto Steve Jobs



Il nemico alle porte

La Consob e la mano invisibile

Altri articoli



PICCOLE GRANDI IMPRESE
DI LUCA PAGNI

La grande sfida del teleshopping

La crisi colpisce anche i porti turistici
ma siamo sicuri che sia un male?

Altri articoli



PERCENTUALMENTE
DI ROSARIA AMATO

La prova del 9

L'export risolve il Pil, ma non le famiglie

Altri articoli

GLI ESPERTI RISPONDONO

CASA

A cura di Antonella Donati



Compenso extra, quando ne ha diritto l'amministratore

Mia moglie ed il fratello sono proprietari di un appartamento in condominio. Allo stato

Il tuo libro arriva dove
hai sempre sognato.

ilmiolibro.it

24ORE AGI

Roma 19:04
ACEA: NEL I SEMESTRE UTILE NETTO +52,1% A 2010, MLN

Parigi 18:42
AIR FRANCE-KLM: TORNA IN UTILE NEL PRIMO TRIMESTRE

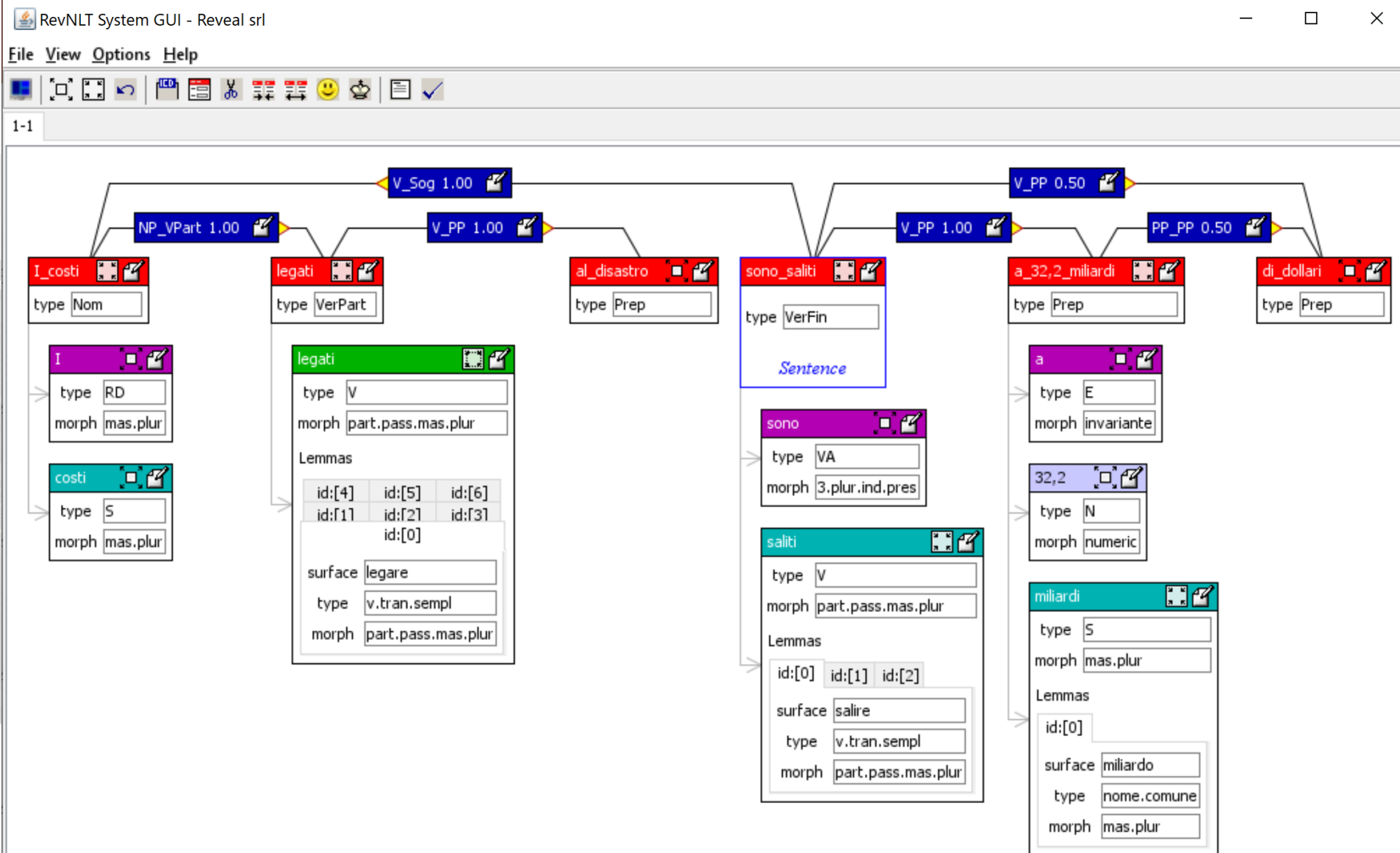
← 3 → Le altre notizie

CREDITO ALLE IMPRESE

Microimprese: con la crisi aumenta il rischio di credito

IN COLLABORAZIONE CON

NLP: parsing Web texts



Course Laboratories:

- During the Course some laboratory sessions (4 hours per week) are scheduled jointly with the *Machine Learning* course (Gambosi):
 - Machine Learning platforms: Pytorch, KERP, SciKit
 - NLP tools: Spacy, RevNLT (Reveal s.r.l.)
 - Deep Learning Tools for NLP:
 - Recursive Neural Networks for the acquisition of Domain specific Dictionaries
 - Transformers for Semantic Parsing and Natural Language Inference (e.g. paraphrasing)
 - Named Entity recognition and Wikification
 - Deep Learning for Visual Recognition
 - Object Detection, Captioning, Visual Question Answering
 - Deep Learning for Web Applications:
 - Sentiment Analysis
 - Fake News Detection

Un esempio: Kelp: Java-based kernel framework

GitHub

Explore Features Enterprise Pricing

Sign up

Sign in

KeLP

KeLP (Kernel-based Learning Platform) is a Java machine learning platform developed within the SAG group and the QCRI.

University of Roma, Tor Vergata <http://sag.art.uniroma2.it/demo-software/kelp>

Repositories

People 3

Filters

kelp-additional-algorithms

Updated 9 days ago

kelp-full

Updated 12 days ago

Semantic Analytics Group @ Uniroma2

SAG is the Semantic Analytics Group at the University of Rome, Tor Vergata

People Research Teaching Publications Projects Demo & Software Contacts

KeLP (Kernel-based Learning Platform)

KeLP (Kernel-based Learning Platform) is a machine learning platform developed within the SAG group. It is entirely written in Java and it is strongly focused on *Kernel Machines*. It includes different Online and Batch Learning and Classification algorithms, Kernel functions, ranging from vector-based to structural kernels. KeLP allows to build complex kernel machine based systems, leveraging on the Java language and on a JSON interface to store and load classifiers configurations as well as to save the models to be reused.

For a deeper look, you can visit [What's inside KeLP page](#).

Downloads

KeLP is released under [Maven](#). To use it, please refer to the [Installation page](#). To download KeLP source code you can go to the github [KeLP page](#).

Authentication

[Log In](#)

News

- SAG's KeLP team ranked first at the [SemEval 2016 Community Question Answering Task](#) February 16, 2016
- [KeLP 2.0.2 released!](#) February 16, 2016
- [KeLP 2.0.1 released](#) January 13, 2016
- [The ECIR 2016 paper has been accepted!](#) December 30, 2015
- [KeLP 2.0.0 released](#) December 4, 2015
- [SAG with Reveal @ Maker Faire 2015, Rome!](#) October 16, 2015

<https://github.com/SAG-KeLP>

<http://sag.art.uniroma2.it/demo-software/kelp/>

KELP applications: cQA

General Description | Subtasks | Data and Tools | Important Dates | **Results** | Call for Papers

SemEval-2016 Task 3

General Description | Subtasks | Data and Tools

SemEval-2016 Task 3

Task 3: Community Question Answering

Building on the success of [SemEval 2015 Task 3](#) "Answer Selection in Community Question Answering" (see [the task description paper](#)), we propose an extension, which covers a full task on Community Question Answering (CQA) and which is, therefore, closer to a real application (see, e.g., [Qatar Living forum](#)).

CQA systems are gaining popularity online. Such systems are seldom moderated, quite open, and thus they have little restrictions, if any, on who can post and who can answer a question. On the positive side, this means that one can freely ask any question and expect some good, honest answers. On the negative side, it takes effort to go through all possible answers and to make sense of them. For example, it is not unusual for a question to have hundreds of answers, which makes it very time-consuming for the user to inspect and to winnow through them all. The present task could help to automate the process of finding good answers to new questions in a community-created discussion forum (e.g., by retrieving similar questions in the forum and by identifying the posts in the answer threads of those similar questions that answer the original question well).

In essence, the main CQA task can be defined as follows:

"given (i) a new question and (ii) a large collection of question-comment threads created by a user community, rank the comments that are most useful for answering the new question"

Results

- └ The evaluation results can be found [here](#)
- └ The gold labels, submissions and scores for all teams can be found [here](#)
- └ The gold labels inside the test XML can be found [here](#)

Task participants are strongly encouraged to submit a system description to the 2016:

<http://alt.qcri.org/semeval2016/index.php?id=call-for-papers>

KELP applications: cQA

[General Description](#)[Subtasks](#)[Data and Tools](#)[Important Dates](#)[Results](#)[Call for Papers](#)

SemEval-2016 Task 3

Results

- └ The evaluation results can be found [here](#)
- └ The gold labels, submissions and scores for all teams can be found [here](#)
- └ The gold labels inside the test XML can be found [here](#)

Task participants are strongly encouraged to submit a system description paper by March 4, 2016:

<http://alt.qcri.org/semEval2016/index.php?id=call-for-papers>

Contact Info

Organizers

- ▶ Preslav Nakov, Qatar Computing Research Institute, HBKU
- ▶ Lluís Màrquez, Qatar Computing Research Institute, HBKU
- ▶ Alessandro Moschitti, Qatar Computing Research Institute, HBKU
- ▶ Walid Magdy, Qatar Computing Research Institute, HBKU
- ▶ James Glass, CSAIL-MIT
- ▶ Bilal Randeree, Qatar Living

email : *semEval-cqa@googlegroups.com*

Other Info

Announcements

KELP

applications: cQA

General Description	Subtasks
SemEval-2016 Task 1	

Team ID	Team Affiliation
ConvKN	Qatar Computing Research Institute,
ECNU	East China Normal University, China
ICL00	Institute of Computational Linguistics
ICRC-HIT	Intelligence Computing Research Center
ITNLP-AiKF	Intelligence Technology and Natural Language Processing
Kelp	University of Roma, Tor Vergata, Italy
MTE-NN	Qatar Computing Research Institute,
overfitting	University of Waterloo, Canada
PMI-cool	Sofia University, Bulgaria
QAIIIIT	IIIT Hyderabad, India
QU-IR	Qatar University, Qatar
RDI.team	RDI Egypt, Cairo University, Egypt
SemanticZ	Sofia University, Bulgaria
SLS	MIT Computer Science and Artificial Intelligence Laboratory
SUPer.team	Sofia University, Bulgaria; Qatar Computing Research Institute
UH-PRHLT	Pattern Recognition and Human Language Understanding
UniMelb	The University of Melbourne, Australia
UPC_USMBA	Universitat Politècnica de Catalunya

Table 5: The participating teams.

Submission	MAP	AvgRec	MRR	P	R	F1	Acc
1 Kelp-primary	79.19₁	88.82₁	86.42₁	76.96₁	55.30₈	64.36₅	75.11₂
ConvKN-contrastive1	78.71	88.98	86.15	77.78	53.72	63.55	74.95
SUPer.team-contrastive1	77.68	88.06	84.76	75.59	55.00	63.68	74.50
2 ConvKN-primary	77.66₂	88.05₃	84.93₄	75.56₂	58.84₆	66.16₂	75.54₁
3 SemanticZ-primary	77.58₃	88.14₂	85.21₂	74.13₄	53.05₁₀	61.84₈	73.39₅
ConvKN-contrastive2	77.29	87.77	85.03	74.74	59.67	66.36	75.41
4 ECNU-primary	77.28₄	87.52₅	84.09₆	70.46₆	63.36₄	66.72₁	74.31₄
SemanticZ-contrastive1	77.16	87.73	84.08	75.29	53.20	62.35	73.88
5 SUPer.team-primary	77.16₅	87.98₄	84.69₅	74.43₃	56.73₇	64.39₄	74.50₃
MTE-NN-contrastive2	76.98	86.98	85.50	58.71	70.28	63.97	67.83
SUPer.team-contrastive2	76.97	87.89	84.58	74.31	56.36	64.10	74.34
MTE-NN-contrastive1	76.86	87.03	84.36	55.84	77.35	64.86	65.93
SLS-contrastive2	76.71	87.17	84.38	59.45	67.95	63.41	68.13
SLS-contrastive1	76.46	87.47	83.27	60.09	69.68	64.53	68.87
6 MTE-NN-primary	76.44₆	86.74₇	84.97₃	56.28₉	76.22₁	64.75₃	66.27₈
7 SLS-primary	76.33₇	87.30₆	82.99₇	60.36₈	67.72₃	63.83₆	68.81₇
ECNU-contrastive2	75.71	86.14	82.53	63.60	66.67	65.10	70.95
SemanticZ-contrastive2	75.41	86.51	82.52	73.19	50.11	59.49	72.26
ICRC-HIT-contrastive1	73.34	84.81	79.65	63.43	69.30	66.24	71.28
8 ITNLP-AiKF-primary	71.52₈	82.67₉	80.26₈	73.18₅	19.71₁₂	31.06₁₂	64.43₉
ECNU-contrastive1	71.34	83.39	78.62	66.95	41.31	51.09	67.86
9 ICRC-HIT-primary	70.90₉	83.36₈	77.38₁₀	62.48₇	62.53₅	62.50₇	69.51₆
10 PMI-cool-primary	68.79₁₀	79.94₁₀	80.00₉	47.81₁₂	70.58₂	57.00₉	56.73₁₂
UH-PRHLT-contrastive1	67.57	79.50	77.08	54.10	50.11	52.03	62.45
11 UH-PRHLT-primary	67.42₁₁	79.38₁₁	76.97₁₁	55.64₁₀	46.80₁₁	50.84₁₁	63.21₁₀
UH-PRHLT-contrastive2	67.33	79.34	76.73	54.97	49.13	51.89	62.97
12 QAIIIIT-primary	62.24₁₂	75.41₁₂	70.58₁₂	50.28₁₁	53.50₉	51.84₁₀	59.60₁₁
QAIIIIT-contrastive2	61.93	75.22	69.95	49.48	49.96	49.72	58.93
QAIIIIT-contrastive1	61.80	75.12	69.76	49.85	50.94	50.39	59.24
Baseline 1 (IR)	59.53	72.60	67.83	—	—	—	—
Baseline 2 (random)	52.80	66.52	58.71	40.56	74.57	52.55	45.26
Baseline 3 (all 'true')	—	—	—	40.64	100.00	57.80	40.64
Baseline 4 (all 'false')	—	—	—	—	—	—	59.36

Table 1: **Subtask A, English (Question-Comment Similarity):** results for all submissions. The first column shows the rank of the primary runs with respect to the official MAP score. The second column contains the team’s name and its submission type (primary vs. contrastive). The following columns show the results for the primary, and then for other, unofficial evaluation measures. The subindices show the rank of the primary runs with respect to the evaluation measure in the respective column.

References

- Mitchell, Tom. M. 1997. *Machine Learning*. New York: McGraw-Hill.
- [Kernel machines, neural networks and graphical models](#), P. Frasconi, A. Sperduti, A. Starita, Rivista AI*IA Numero speciale per i “50 anni di IA”, 2007.
- Very good video lectures by Andrew Ng (Stanford) <http://academicearth.org/courses/machine-learning>