

Web Mining

Introduction and Course Outline

Simone Paolo Ponzetto Anne Lauscher Chia-Chien Hung

FSS 2020

Since 2016 full professor of Information Systems (Wifo) in Mannheim

- Main research areas
 - Knowledge Acquisition
 - Natural Language Processing
 - Computational Social Science



- TAs: Anne Lauscher, Chia-Chien Hung
- {anne, chia-chien, simone}@informatik.uni-mannheim.de







Introduction and Course Organization

- **1. Course Organization**
- 2. The World Wide Web
 - **1.** The Classic Document Web
 - 2. The Web of Data
 - **3.** Web 2.0 Applications
- 3. What is Web Mining?
 - 1. Web Usage Mining
 - 2. Web Structure Mining
 - **3. Web Content Mining**

Lecture (IE 671)

- covers different types of Web Mining methods
- presents examples of Web Mining applications
- discusses how to evaluate learned models

Labs

students get their hand dirty with tools and code

Evaluation

60 min final exam

Project (IE 684)

- teams of four students realize a Web Mining project
- teams may choose their own data sets and tasks (in addition, we will propose some suitable data sets and tasks)
- write a summary about the project, present the project results

Evaluation

Implementation + report + presentation

Materials

ILIAS eLearning System, https://ilias.uni-mannheim.de/

Time and Location

Lecture:

Tuesday, 10:15 to 11:45, Room: B 6, A104

Labs:

Tuesday, 15:30 to 17:00, Room: B 6, A104

ILIAS					
UNIVERSITÄT MANNHEIM	PERSÖNLICHER SCI	HREIBTISCH 👻	PORTAL ²	MAGAZIN 👻	HILFE
Übersicht					
Neuigkeiten - Letzte 5 Tage	•	Meine Mitglie	edschaften		
0 Neuigkeit(en)		Fachbereich Bet	riebswirtschaftsl	ehre	

Schedule

Morning session (10:15-11:45)		Afternoon session (15:30-15:00)	
11/2/20	Lecture: Introduction to Web Mining	11/2/20	-
18/2/20	Lecture: Recommender Systems	18/2/20	Lab: Recommender Systems
25/2/20	Lecture: Recommender Systems	25/2/20	Lab: Recommender Systems
3/3/20	Lecture: Social Network Analysis	3/3/20	Lab: Social Network Analysis
10/3/20	Lecture: Web Content Mining	10/3/20	Lab: Web Content Mining
17/3/20	Lecture: Web Content Mining	17/3/20	Lab: Web Content Mining
24/3/20	Lecture: Web Content Mining	24/3/20	Lab: Web Content Mining
31/3/20	Introduction to Student Projects		-
7/4/20		Osterferien	
14/4/20		Osterferien	
21/4/20	Feedback about Project Oulines	21/4/20	-
28/4/20	-	28/4/20	-
5/5/20	Coaching session (optional)	5/5/20	
12/5/20	Coaching session (optional)	12/5/20	
19/5/20	Coaching session (optional)	19/5/20	
26/5/20	Presentation of the projects	26/5/20	Presentation of the projects

Lecture Videos

Video recordings of past lectures (outdated but useful as extra complementary materials)

http://dws.informatik.uni-mannheim.de/en/teaching/lecturevideos/



Software/libraries

You will not succeed in this module if you do no know how to program (we assume some fluency in Python)...



surprise

A Python scikit for recommender systems

Home Documentation

⊖ GitHub page ★ Star ¥Fork

Maintained by Nicolas Hug Page built with Jekyll and H



Making neural nets uncool again

<u>Home</u> <u>About</u> <u>Our MOOC</u> <u>Posts by Topic</u>

© fast.ai 2020. All rights reserved.



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable BSD license

Questions?



The Web is a global information space build on a set of technical standards for the identification, retrieval and representation of content.



- Hypertext Transfer Protokoll (HTTP): Protocol for interacting with Web resources.
- **Content Formats: HTML, XML, RDF, ...**
- The Web was invented in 1989 at CERN by Tim Berners-Lee



Architectural Principles of the Web

Architecture of the World Wide Web, Volume One, W3C Specification, 2004, http://www.w3.org/TR/webarch/

Topology of the Web Today



2.1 The Classic Document Web

Global information space consisting of interlinked Web documents (text, images, multimedia).



Size of the Indexed Web: approx. 60 billion pages

http://www.worldwidewebsize.com/

- Estimated on the basis of a method that combines word frequencies from a corpus and search counts returned by the engines
- **50** words are sent to all four search engines
- The number of webpages found for these words are recorded
- Multiple extrapolated estimations are made of the size of the engine's index, which are subsequently averaged
- Example: The word "the" appears in 67% of all English pages and has 25.2 billion hits on Google

Link Structure of the Web: In-Degree

The link distribution follows (kind of) a power law.

- A small number of pages is target of many links.
- A large number of pages is target of only a few or no links.

Classic Paper:

- Broder at al.: Graph Structure in the Web. WWW2000.
- AltaVista crawl with over 200 million pages and 1.5 billion links
- Conclusion: Log-log scale plot shows power-law.



Broder et al. (2000)

Power law with exponent 2.1

(200 million pages and 1.5 billion links from Altavista crawl 2000) WDC Hyperlink Graph (2012) Best power law exponent 2.24 (3 billion pages and 128 billion links from Common Crawl 2012)



Link Structure of the Web: Bow-Tie

Four mayor components (Border at al., WWW2000)

Central Strongly Connected Component (SCC)

- pages that can reach one another along directed links
- about 30% of the Web (normal pages)
- IN Group
 - can reach SCC but cannot be reached from it
 - about 20% (maybe new pages or boring ones)

OUT Group

- can be reached from SCC but cannot reach it
- about 20% (maybe company pages that don't link)

Tendrils

- cannot reach SCC and cannot be reached by it
- about 20%
- Unconnected
 - about 10%



Probability of path between nodes is 24%

Web Graph with Strongly Connected Components



A strongly connected component (SCC) in a directed graph is a subset of the nodes such that:

- 1. every node in the subset has a path to every other node
- 2. the subset is not part of some larger set with the property that every node can reach every other.

Largest Strongly Connected Component



From a Web of Documents to a Web of Data

Web of documents



"Documents"

Key characteristics:

- 1. <u>Names</u> (URIs)
- 2. <u>Documents</u> (Resources) described by HTML, XML, etc.
- 3. Interactions via HTTP
- 4. <u>(Hyper)Links</u> between documents or anchors in these documents

Web of Documents vs. Web of Data



Web of data

Web of data

Key characteristics:

- Links between arbitrary things (e.g., persons, locations, events, buildings)
- Structure of data on Web pages is made explicit
- Things described on Web pages are named and get URIs
- Links between things are made explicit and are typed



2.2. The Web of Data

Available approaches

- 1. <u>semantically markup</u> the content of their HTML pages
- 2. publish structured data in addition to HTML pages



Microformats

Microformat effort dates back to 2003

Small set of fixed formats

- hcard : people, companies, organizations, and places
- XFN : relationships between people
- hCalendar : calendaring and events
- hListing : small-ads; classifieds
- hReview : reviews of products, businesses, events

Key idea

use existing HTML attributes to embed structured data types

Shortcoming

Fixed formats means that it can not represent any kind of data



hCard is a simple, open format for publishing people, companies, organizations on the web, using a 1:1 representation of vCard

```
<div class="vcard">
  <a class="fn org url" href="http://www.commerce.net/">CommerceNet</a>
  <div class="adr">
    <span class="type">Work</span>:
    <div class="street-address">169 University Avenue</div>
    <span class="locality">Palo Alto</span>,
    <abbr class="region" title="California">CA</abbr>&nbsp;&nbsp;</a>
    <span class="postal-code">94301</span>
    <div class="country-name">USA</div>
  </div>
  <div class="tel">
  <span class="type">Work</span> +1-650-289-4040
  </div>
  <div class="tel">
    <span class="type">Fax</span> +1-650-289-4041
  </div>
  <div>Email:
   <span class="email">info@commerce.net</span>
  </div>
</div>
```

RDFa

serialization format for embedding RDF data into HTML pages



- proposed in 2004, W3C Recommendation in 2008
- can be used together with any vocabulary

```
1 <html xmlns="http://www.w3.org/1999/xhtml"
2 xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
3 xmlns:foaf="http://xmlns.com/foaf/0.1/">
4 ...
5 <div about="http://example.com/Peter" typeof="foaf:Person">
6 <span property="foaf:name">Peter Smith</span> knows
7 <span property="foaf:name">Peter Smith</span> knows
7 <a rel="foaf:knows" href="http://example.com/Paula">Paula
Jones</a>.
8 <//div>
9 ...
```

- alternative technique for embedding structured data
- proposed in 2009 by WHATWG as part of HTML5 work
- tries to be simpler than RDFa (5 new attributes instead of 8)

HTML

Schema.org

- Problem: which vocabulary to use?
- **Schema.org** provides a collection of shared vocabularies
- Launched in June 2011 by Bing, Google and Yahoo

- Create a common set of schemas for webmasters to mark-up with structured data their websites.
- 200+ Types: Event, Organization, Person, Place, Product, Review
- Encoding: Microdata, RDFa, JSON-LD

Schema.org: example schema

schema.org	Custom Sea	rch Q
Home	Schemas	Documentation

Person

Thing > Person

A person (alive, dead, undead, or fictional).

[more...]

Property	Expected Type	Description
Properties from Person		
additionalName	Text	An additional name for a Person, can be used for a middle name.
address	PostalAddress or Text	Physical address of the item.
affiliation	Organization	An organization that this person is affiliated with. For example, a school/university, a club, or a team.
alumniOf	EducationalOrganization or Organization	An organization that the person is an alumni of. Inverse property: <u>alumni</u> .
award	Text	An award won by or for this item. Supersedes awards.
birthDate	Date	Date of birth.
birthPlace	Place	The place where the person was born.

Schema.org: examples

RDFa

Microdata

Usage of Schema.org Data @ Google

Gramercy Tavern - Flatiron - New York, NY | Yelp

www.yelp.com > Restaurants > American (New) ▼ ★★★★★ Rating: 4.5 - 1,288 reviews - Price range: \$\$\$\$ Jeff C and I were in New York for vacation, and I wanted to treat him to a nice dinner for Gramercy Tavern is certainly a legendary NY dining establishment.

Gramercy Tavern Restaurant - New York, NY | OpenTable

www.opentable.com > … > Gramercy restaurants ▼ ★★★★ Rating: 4.7 - 508 reviews - Price range: \$50 and over Book now at Gramercy Tavern in New York, explore menu, see photos and read 508 reviews: "The menu was so limited but it was worth trying, food was deli..."

Data snippets within search results

Data snippets within info boxes



The Black Keys

Band

The Black Keys is an American rock duo formed in Akron, Ohio in 2001. The group consists of Dan Auerbach and Patrick Carney. Wikipedia

Origin: Akron, Ohio, United States

Members: Dan Auerbach, Patrick Carney

Record labels: Fat Possum Records, Nonesuch Records, V2 Records, Alive Naturalsound Records

Awards: Grammy Award for Best Rock Album, more

Upcoming events

Jun 20 Fri	The Black Keys Neuhausen ob Eck (near you)	
May 16 ^{Fri}	The Black Keys Gulf Shores, AL	
Jun 22 Sun	The Black Keys Scheeßel	

Extend the Web with a single global data graph

- 1. by using RDF to publish structured data on the Web
- 2. by setting links between data items within different data sources.



The Web Today as a Multitude of Data Silos



From Data Silos vs. the Web of Data

- 1. <u>Use URIs as names</u> for things.
- 2. <u>Use HTTP URIs</u> so that people can look up those names.
- 3. When someone looks up a URI, provide useful RDF information.
- 4. Include RDF <u>statements that link to other URIs</u> so that they can discover related things.



Entities are identified with HTTP URIs



HTTP URIs take the role of global primary keys.

pd:cygri = http://richard.cyganiak.de/foaf.rdf#cygri dbpedia:Berlin = http://dbpedia.org/resource/Berlin

URIs can be looked up on the Web



By following RDF links applications can

navigate the global data graph

URIs can be linked to navigate the data space



By following RDF links applications can

- navigate the global data graph
- discover new data sources

2.3. Web 2.0 Applications

- A multitude of Webbased applications has sprung up which enable users to share information.
- These applications
 - collect large amounts of data using proprietary schemata.
 - form separate data spaces that are only partly accessible from the Web via:
 - **1. HTML interfaces**
 - 2. Web APIs



HTML Interfaces

Allows browser-based access to profile, communication, etc.

y			Home				*\$	Q Search Twitter	
٥	Home			What's happe	ening?			Trends for you	Ô
#	Explore				٢		et	Trending in Germany #AKKRücktritt 1,587 Tweets	
Q	Notification	S	ş	Frankfurter Allg	gemeine 🔗 @faznet • 1m	·		Michael Piotrowski is Tweeting about this	
\square	Messages		c	Miese Zahlen, m schwedische Vo	aue Aussichten: Bei Daim rstandschef Ola Källenius nger ist nicht ganz unsch	ler läuft nicht alles rund. Der hat weiter schwere Tage vor uldig, schreibt @MeckGeorg:		Trending in Germany Weltuntergang	
	Bookmark	T		Following				Trending in Germany #Sabine 56 9K Tweats	
Ē	Lists	Frankfurter A @faznet Die wichtigste	llgemeine n Nachricl	e 🍼				 ZEIT ONLINE and tagesschau are Tweeting about this 	
	Profile	Tages, die bes Empfehlungen	ten Faz.ne der Reda	et-Artikel und ktion.	0			Politics · Trending #Merz	
	More	Impressum: fa Datenschutz: f	z.net/impr faz.net/da	essum tenschutz				2,267 Tweets left Postillon is Tweeting about this	
	Tweet	296 Following Followe Lauterb	539.61 d by Kristia ach, and 4	Followers n Kersting, Karl others you follow	ichten: Daimler im Abstie	rgskampf		Celebrities · Trending Joaquin Phoenix 501K Tweets	
	Theor			Miese Zahlen, schwedische V & faz.net	maue Aussichten: Bei Dai ⁄orstandschef Ola Källeniu	mler läuft nicht alles rund. De Is hat weiter schwere Tage vo	er	Show more	

Web APIs

Access the data programmatically



Web APIs slice the Web into Data Silos



Definition

Acquiring useful information from

- Web content
- Web structure
- Web usage data

Web data pose unique challenges:

- 1. Large volume
- 2. <u>Semi-structured</u>
- 3. <u>Heterogeneous</u>
- 4. Distributed

Web Mining is a Multi-Disciplinary Field

Draws ideas and techniques from

Sub-Fields

- 1. Web Usage Mining
- 2. Web Structure Mining
- 3. Web Content Mining



Definition

Discovery of patterns in data collected or generated as a result of user interactions with one or more web sites.

Sources of Data

- **1.** automatically generated data stored in server access logs
- 2. e-commerce and product-oriented user events (e.g., shopping cart changes, ad or product click-throughs, purchases)
- **3.** user profiles (e.g. Facebook) and/or user ratings (likes)
- 4. page attributes, page content, site structure
- 5. additional domain knowledge and demographic data

Leading Usage Data Collections







Enable the

- analysis of the current interests and behavior of the world's population.
- identification of suspected terrorists.

The Web Usage Mining Process



Example: product recommendation



Example: product recommendation

• • •	< > Q Search			ඹ spponzetto 🗸 🗸
 G Home O Browse (∞) Radio YOUR LIBRARY Made For You Recently Played 	Your Made for spponzetto Discover Our weekly mixtape of fresh Made for spponzetto by Spoti Made for spponzetto by Spoti	Weekly music. Enjoy new discoveries and de s! fy • 30 songs, 2 hr 3 min	eep cuts chosen just for you. Up	dated every FOLLOWER 1
Liked Songs				
Artists	Q Filter			Download
Podcasts	TITLE	ARTIST	ALBUM	
	Dancing in the Street - 2002 Remaster	David Bowie, Mick Jagger	Dancing In The Street E.P.	15 hours ago
PLAYLISTS Anime&Games ►	♡ Eloise	Tino Casal	Grandes Exitos (Etiqueta N	15 hours ago
Artists 🕨	♡ Theme From Starsky & Hutch - Funky People Mix	James Taylor Quartet	Wait A Minute	15 hours ago
ввс 🕨	Pop Muzik - Nik Launay '79 12"	M, Robin Scott	New York, London, Paris, M	15 hours ago
Books / Movies	♡ Get Up (I Feel Like Being A) Sex Machine - Pt. 1 & 2	James Brown	20 All-Time Greatest Hits!	15 hours ago
Classics	\heartsuit Indiana Jones and the Temple of Doom	Royal Philharmonic Orchestra	Movie Legends: The Music	15 hours ago
Discover 👻	♡ Juliet	Robin Gibb	Love Songs	15 hours ago
Discover Weekly	🌣 🛛 Just a Gigolo / I Ain't Got Nobody - 45 Version	David Lee Roth	Just A Gigolo/I Ain't Got N	15 hours ago
Ninja Tune - On	♡ Gold	Spandau Ballet	Gold - The Best of Spandau	15 hours ago
Warp Selections	♡ Jerusalem	Fat Les	Jerusalem	15 hours ago
Funk & Soul 🕨	Don't Stop Believin'	Journey	Escape	15 hours ago
Gilles •	Magnetic Fields, Pt. 2	Jean-Michel Jarre	Les Chants Magnétiques /	15 hours ago
Jazz 🕨	♡ Fiesta	The Pogues	If I Should Fall From Grace	15 hours ago
🕂 New Playlist	♡ My Sharona	The Knack	Get The Knack	15 hours ago

Example: personalized search

Google

🔍 All 🖾 Images 🕞 Videos 🛇 Maps 🖽 News 🗄 More Settings Tools

About 51.700.000 results (0,58 seconds)

Ad · www.confluent.io/ -

kafka

Confluent | Download Apache Kafka® Today | confluent.io

Monitor & Manage Data in Real-Time Using Our Scalable, Reliable, & Flexible Platform. Confluent Platform Includes the Latest Version of Apache Kafka, Plus Enterprise Features. Reduce Ops Burden. Streaming Data Service. Deployable in Minutes. 24/7 Support.

What is Kafka?

More Info. About Kafka & Confluent Platform.

Kafka Definitive Guide

Q

Learn All About Kafka From its Original Developers in this eBook.

kafka.apache.org -

Apache Kafka

Kafka® is used for building real-time data pipelines and streaming apps. It is horizontally scalable, fault-tolerant, wicked fast, and runs in production in thousands ... Introduction · Kafka Streams · Documentation · Apache Kafka Documentation

kafka.apache.org > intro 🔻

Introduction - Apache Kafka - Apache Software

In Kafka the communication between the clients and the servers is done with a simple, highperformance, language agnostic TCP protocol. This protocol is ...

People also ask	
What is Kafka used for?	~
What is meant by Kafka?	~
What is Kafka and how it works?	~
Is Kafka free?	~



....

Apache Kafka is an open-source stream-processing software platform developed by LinkedIn and donated to the Apache Software Foundation, written in Scala and Java. The project aims to provide a unified, high-throughput, low-latency platform for handling real-time data feeds. Wikipedia

License: Apache License 2.0

Developer(s): Apache Software Foundation

Initial release date: January 2011

Stable release: 2.4.0 / December 16, 2019; 48 days ago

Written in: Scala, Java



Example: personalized search

Google	kafka							٩
	Q Alle	🖬 Bilder	▶ Videos	🗉 News	Bücher	: Mehr	Einstellungen	Suchfilter
	Ungefäh	r 50.100.000) Ergebnisse (0,66 Sekund	en)			



de.wikipedia.org > wiki > Franz_Kafka 🔻

Franz Kafka – Wikipedia

Franz Kafka (tschechisch gelegentlich František Kafka, jüdischer Name: אנשיל Anschel; * 3. Juli 1883 in Prag, Österreich-Ungarn; † 3. Juni 1924 in Kierling, ... Hermann Kafka · Kategorie:Franz Kafka · Kafkaesk · Die Verwandlung

kafka.apache.org v Diese Seite übersetzen

Apache Kafka

Kafka® is used for building real-time data pipelines and streaming apps. It is horizontally scalable, fault-tolerant, wicked fast, and runs in production in thousands ... Introduction · Kafka Streams · Apache Kafka Documentation · Documentation

www.franzkafka.de 🔻

Franz Kafka

Dieses Portal zu Franz Kafka bietet zuverlässig und kurzweilig die wichtigsten Informationen über sein Werk.

www.franzkafka.de > franzkafka > das_leben 🔻

Das Leben - Franz Kafka

Juli wird Franz Kafka in Prag geboren. Er ist das erste Kind von Hermann Kafka (1852-1931) und seiner Frau Julie, geb. Löwy (1856-1934). Die jüdischen Eltern ...

www.inhaltsangabe.de > autoren > kafka 🔻

Franz Kafka - Biografie und Inhaltsangaben - Inhaltsangabe.de



Franz Kafka

Schriftsteller

Franz Kafka war ein deutschsprachiger Schriftsteller. Sein Hauptwerk bilden neben drei Romanfragmenten zahlreiche Erzählungen. Wikipedia

Geboren: 3. Juli 1883, Prag, Tschechien

Gestorben: 3. Juni 1924, Kierling, Klosterneuburg, Österreich

Bestattet: 11. Juni 1924, Neuer jüdischer Friedhof, Prag, Tschechien

Kurzgeschichten: Die Verwandlung, Das Urteil, Vor dem Gesetz, MEHR

Beeinflusst von: Fjodor Michailowitsch Dostojewski, MEHR

Bücher

Über 45 weitere ansehen

....

<

Anmelden



3.2 Web Structure Mining

Definition

Discovery of patterns in

- the hyperlink structure of webpages
- the structure of communities that interact on the Web
- Exploits the graph structure, but can of course also be combined with content or usage mining techniques.
- Typical Sources of Data
 - **1.** Web crawls including HTML pages and hyperlinks
 - 2. crawls of the blogosphere
 - 3. social networks including explicit relations between actors (your Facebook friend network)
 - 4. other types of community data (discussion forums, email conversations, ...)

Identification of Prominent Nodes

Question: Who are the "most important" actors in a social network?

Centrality

- A central actor is one involved in many edges.
- The direction of lines is not considered.

Prestige

- A prestigious actor is one who is the target of many arcs.
- The direction of arcs is considered.



A community is a set of actors between which interactions are (relatively) frequent.

Finding a community in a social network is to identify a set of nodes such that they interact with each other more frequently than with those nodes outside the group.

$$9 \xrightarrow{7}{8} \xrightarrow{6}{6} \xrightarrow{2}{2} \xrightarrow{3}{7} \xrightarrow{5}{4} \xrightarrow{3}{4} \xrightarrow{3}{2}$$

- Methods: Components, K-Cores, Islands, …
- Applications: Recommendation based on communities, visualization of huge networks, network compression



Link Prediction

Question: Given a snapshot of a social network, can we infer which new interactions among its members are likely to occur in the near future? (Liben-Nowell & Kleinberg, 2007)

Applications

- Facebook: recommending possible friends
- Tinder: recommending potential matches



Time T Time T+1 D. Liben-Nowell, D. and J. Kleinberg, *The link-prediction problem for social networks.* Journal of the American Society for Information Science and Technology, 58(7) 1019–1031 (2007)

Definition

Automatic extraction of useful information (facts, patterns) from Web content (text, images, multimedia).

Content Mining Tasks

- Content classification and clustering on Web content
- Applications of NLP techniques to social network content

Content Classification

- Supervised Learning: Given a collection of labeled documents/images (training set) find a model for the class as a function of the values of the features.
- Goal: Previously unseen documents/images should be assigned a class as accurately as possible.

Applications

- News categorization
- Topic classification
- Spam detection
- Product categorization

Classification methods commonly used for

Naive Bayes, Support Vector Machines, Deep Neural Nets

Content Clustering

- Unsupervised Learning: Given a set of documents and a similarity measure among documents find clusters such that:
 - documents in one cluster are more similar to one another
 - documents in separate clusters are less similar to one another

Applications

- Topic discovery
- Search result clustering

Techniques

- Algorithms: K-Means, Expectation-Maximization (EM)
- Similarity measures: Cosine, Jaccard

Mixture of Document Clustering and Classification



Sentiment Analysis

The basic task in sentiment analysis is classifying the polarity of a given text at the document, sentence, or feature/aspect level.

Polarity Values

- Positive, neutral, negative
- Stars

Applications

- Document-level: poll prediction from tweets
- Feature/Aspect-level: analysis of product reviews



 Revolution In The Head: The Beatles Records and the Sixties (Englisch) Taschenbuch – 4.

 Dezember 2008

 von Ian MacDonald ~ (Autor)

 ★★★★☆ ~ 175 Sternebewertungen

 • Alle 3 Formate und Ausgaben anzeigen

 Taschenbuch 5,69 €

18 neu ab 5,69 € 5 gebraucht ab 3,60 €

As dazzling as the decade they dominated, The Beatles almost single-handedly created pop music as we know it. Today, their songs are cited as seminal influences by stars like Oasis and Blur. Eloquently giving voice to their time, The Beatles quite simply changed the world.



Topic tracking / hate speech / offensive language

Applications of NLP to social media data

Twitter	% posts	Whisper	% posts
I hate	70.5	I hate	66.4
I can't stand	7.7	I don't like	9.1
I don't like	7.2	I can't stand	7.4
I really hate	4.9	I really hate	3.1
I fucking hate	1.8	I fucking hate	3.0
I'm sick of	0.8	I'm sick of	1.4
I cannot stand	0.7	I'm so sick of	1.0
I fuckin hate	0.6	I just hate	0.9
I just hate	0.6	I really don't like	0.8
I'm so sick of	0.6	I secretly hate	0.7

Twitter		Whisper			
Hate target	% posts	Hate target	% posts		
Nigga	31.11	Black people	10.10		
White people	9.76	Fake people	9.77		
Fake people	5.07	Fat people	8.46		
Black people	4.91	Stupid people	7.84		
Stupid people	2.62	Gay people	7.06		
Rude people	2.60	White people	5.62		
Negative people	2.53	Racist people	3.35		
Ignorant people	2.13	Ignorant peo-	3.10		
		ple			
Nigger	1.84	Rude people	2.45		
Ungrateful people	1.80	Old people	2.18		

Trends for you	Ø
Politics · Trending #Merkel 13.5K Tweets ③ Frankfurter Allgemeine is Tweeting about this	 ✓ S
Politics · Trending #AKKRuecktritt 9,327 Tweets tagesschau, Frankfurter Allgemeine, and 1 m are Tweeting about this	✓ore
Celebrities · Trending Parasite 2.12M Tweets	~
Politics · Trending #BlackRock 15.7K Tweets	~
Trending in Germany #Merz 6,429 Tweets Der Postillon is Tweeting about this	~
Show more	

Questions?

