Apache Mahout Making data analysis easy





Isabel Drost

Nighttime:

Co-Founder, committer Apache Mahout. Organiser of Berlin Hadoop Get Together.

Daytime:

Software developer. Guest lecturer at TU Berlin. Co-Organiser Berlin Buzzwords 2010.









Machine learning background?



Agenda

• Data Mining/ Machine Learning?

• Why is scaling hard?

• Introducing Apache Hadoop.

• Going beyond simple statistics.

Machine learning – what's that?

Data Mining Applications

- Marketing.
- Surveillance.
- Fraud Detection.
- Scientific Discovery.
- Discover items usually purchased together.

= Extracting patterns from data.

Machine Learning Applications

- E-Mail spam classification.
- News-topic discovery.
- Building recommender systems.

= Extracting prediction models from data.



Image by John Leech, from: The Comic History of Rome by Gilbert Abbott A Beckett. Bradbury, Evans & Co, London, 1850s Archimedes taking a Warm Bath

Archimedes model of nature

 $\frac{Density of Object}{Density of Fluid} = .$

Weight Weight – Apparent immersed weight



Nog

Wille.

Doy u.

1.61.

car por

m.fot.

de ann Highe

11. cur poll .cc

tt . soc. fet.

sm/ Suge.

A HIN







18



June 25, 2008 by chase-me A STREET http://www.flickr.com/photos/sasy/2609508999



An SVM's model of nature



The challenge

Mission

Provide scalable data mining algorithms.



January 8, 2008 by Pink Sherbet Photography http://www.flickr.com/photos/pinksherbet/2177961471/

Illin

Massive data as in:

Cannot be stored on single machine. Takes too long to process in serial.

Idea: Use multiple machines.

Challenges when scaling out.

................ 1111 1..... 111111111111 ngle mach ines tend to fail PowerEdge 111111 Hard 2650 di S 11111 er D S

1.

100

....

@ po 1650

...

........

m

More machines – increased failure probability.

January 11, 2007, skreuzer http://www.flickr.com/photos/skreuzer/354316053/

10

63

Typical developer



- Has never dealt with large (petabytes) amount of data.
- Has no thorough understanding of parallel programming.
- Has no time to make software production ready.

http://www.flickr.com/photos/jaaronfarr/3384940437/ March 25, 2009 by jaaron

February 29, 2008 by Thomas Claveirole http://www.flickr.com/photos/thomasclaveirole/2300932656/

http://www.flickr.com/photos /jaaronfarr/3385756482/ March 25, 2009 by jaaron

May 1, 2007 by danny angus http://www.flickr.com/photos/killerbees/479864437/



http://www.flickr.com/photos/cspowers/282944734/ by cspowers on October 29, 2006

Easy distributed programming.

Well known in industry and research.

Scales well beyond 1000 nodes.



Hadoop assumptions



Moving computation is cheap. Moving data is expensive.



Ideas:

Move computation to data. Write software that is easy to distribute.

Assumptions:

Systems run on spinning hard disks. Disk seek >> disk scan.



Ideas:

Improve support for large files. File system API makes scanning easy.





HDFS building blocks







(Graphics: Thanks to Thilo.)

Anatomy of a file write



Slide inspired by: "Hadoop – The definitive guide", Tom White, O'Reilly

Anatomy of a file write



Slide inspired by: "Hadoop – The definitive guide", Tom White, O'Reilly

Anatomy of a file write



Slide inspired by: "Hadoop – The definitive guide", Tom White, O'Reilly
Anatomy of a file write



Slide inspired by: "Hadoop – The definitive guide", Tom White, O'Reilly

Map/Reduce by example



?xml version="1.0" encoding="UTF-8"?

<copml version="1.0" >

<head>

<text></text>

</head>

⊲body>

dutline htmlUrl="http://eventseer.net" title="EventSeer - A Digital Library of Call for Papers" useC alDefault" version="RSS" type="rss" xmlUrl="http://eventseer.net/feeds/main/rss.xml" id="312053548" tex tseer.net" />

dutline isOpen="false" id="669809145" text="Silent" >

<outline htmlUrl="http://www.theserverside.com" title="TheServerSide.com: Patterns" useCustomFetchIn ersion="RSS" type="rss" xmlUrl="http://www.theserverside.com/rss/theserverside-j2eepatterns-rss2.xml" i taining up-to-date news, discussions, patterns, resources, and media" />

doutline htmlUrl="http://chadwa.wordpress.com" title="Chad's Search Blog" useCustomFetchInterval="fa
S" type="rss" xmlUrl="http://chadwa.wordpress.com/feed/" id="545368194" text="Chad's Search Blog" descr
" />

dutline htmlUrl="http://www.find23.net/Site/Blog/Blog.html" title="My Blog" useCustomFetchInterval=
"RSS" type="rss" xmlUrl="http://www.find23.net/Site/Blog/rss.xml" id="1620106192" text="My Blog" description")

doutline htmlUrl="http://emotion.inrialpes.fr/~dangauthier/blog" title="Yet Another Machine Learning
eMode="globalDefault" version="RSS" type="rss" xmlUrl="http://emotion.inrialpes.fr/~dangauthier/blog/fe
g" />

dutline htmlUrl="http://ml.typepad.com/machine_learning_thoughts/" title="Machine Learning Thoughts ="globalDefault" version="RSS" type="rss" xmlUrl="http://ml.typepad.com/machine_learning_thoughts/rss.xmlurl="http://ml.typepad.com/machine_learning_thoughts/rss.xmlurl="http://ml.typepad.com/machine_learning_thoughts/rss.xmlurl="http://ml.typepad.com/machine_learning_thoughts/rss.xmlurl="http://ml.typepad.com/machine_learning_thoughts/" title="Machine_learning_thoughts/rss.xmlurl="http://ml.typepad.com/machine_learning_thoughts/rss.xmlu

doutline htmlUrl="http://yaroslavvb.blogspot.com/" title="Machine Learning, etc" useCustomFetchInter ion="RSS" type="rss" xmlUrl="http://yaroslavvb.blogspot.com/feeds/posts/default" id="805998569" text="M doutline htmlUrl="http://ptufts.blogspot.com/" title="Pinhead's Progress" useCustomFetchInterval="fa S" type="rss" xmlUrl="http://ptufts.blogspot.com/" title="Misc Research Stuff" useCustomFetchInterval="fa soutline htmlUrl="http://resnotebook.blogspot.com/" title="Misc Research Stuff" useCustomFetchInterval on="RSS" type="rss" xmlUrl="http://resnotebook.blogspot.com/" title="Misc Research Stuff" useCustomFetchInterval on="RSS" type="rss" xmlUrl="http://resnotebook.blogspot.com/" title="Absolutely Regular" useCustomFetchInterval doutline htmlUrl="http://absolutely-regular.blogspot.com/" title="Absolutely Regular" useCustomFetch version="RSS" type="rss" xmlUrl="http://absolutely-regular.blogspot.com/" title="Absolutely Regular" useCustomFetch version="RSS" type="rss" xmlUrl="http://absolutely-regular.blogspot.com/" title="Absolutely Regular" useCustomFetch version="RSS" type="rss" xmlUrl="http://absolutely-regular.blogspot.com/feeds/posts/default" id="17850! doutline htmlUrl="http://atomai.blogspot.com/" title="Data Mining, Analytics and Artificial Intellige Mode="globalDefault" version="RSS" type="rss" xmlUrl="http://atomai.blogspot.com/feeds/posts/default" in nt in data mining, artificial intelligence, analytics, intelligent agents, semiconductors, distributing siness Objects, Oracle, Intel, AMD, or Pentaho. Heuristic, Six Sigma, or CMM. Contractor or in-house. H ail_com" /> isabel@h1349259:~\$ more data/feeds.opml | grep -o "http://[0-9A-Za-z\-_\.]*" | s

- ort | uniq --count | sort | tail
 - 3 http://agbs.kyb.tuebingen.mpg.de
 - 3 http://irgupf.com
 - 3 http://jeffsutherland.com
 - 4 http://ml.typepad.com
 - 4 http://weblogs.java.net
 - 4 http://www.gridvm.org
 - 4 http://yaroslavvb.blogspot.com
 - 5 http://feeds.feedburner.com
 - 6 http://blogsearch.google.com
 - 10 http://arxiv.org

pattern="http://[0-9A-Za-z\-_\.]*"

grep -o "\$pattern" feeds.opml | sort | uniq --count





grep -o "\$pattern" feeds.opml	sort	uniqcount
ΜΑΡ	SHUFFLE	R E D U C E

pattern="http://[0-9A-Za-z\-_\.]*"



Mission

Provide scalable data mining algorithms.



HowTo: From data to information.

January 3, 2006 by Matt Callow http://www.flickr.com/photos/blackcustard/81680010

COMMUNITY NEWS

Finishing touches still to come

A glimpse of today, yesterday

http://www.flickr.com/photos/redux/409356158/









(in Datassuche Immobilien Autometici Joka Beineurgeboie

0

ZEIT CONLINE DATENSCHUTZ

<text><text><text><text> heckpoints and current best









Annaldes | Basistria erschicken | Facebook, Twitten

Development and the series of the series of the series and the series of understanding of the rationale underpinning heckpoints and current best









Finally, there h lernet themselve agents are a and how to Tausende demon Bürgerrechte im I particular Für einen besseren Arbeit larger tex skarte: 130 hey sho

Rund 7500 Demonstranten nahmen an dem Prot "Freiheit statt Angat – Stoppt den Überwachungs

experteer^{de}

reiheit stati

Rund 7500 Demonstranten nahmen an dem Protestrug unter dem Mo "Freiheit statt Angst – Stoppt den Überwachungswahn" in Berlin teil.

experteerde 🛄

Merivale mu By JUSTIN S. CAMPBELL Error of the second

ZEIT DATENSCHUTZ

Partnersuche

DATU!

E. QUELL

ROMM

* EMPFE

Buzz

Datens

NEU IM R

im Neta

2. FUTUR URHEB

gut"

. SPAM

den Hå

IPHON

NEU AUF

I. NACHR

STARTSEITE POLITIK WIRTSCHAFT MEINUNG GESELLSCHAFT KULTUR WISSEN DIGITAL STUDIL SPORT

Datenschutz Games Internet Mobil

DATENSCHUTZ

Tausende demonstrieren für Bürgerrechte im Netz

Für einen besseren Arbeitnehmerdatenschutz und gegen die Gesundheitskarte: 130 Organisationen hatten zur SCHLA Demonstration aufgerufen. Sie fürchten den Überwachungsstaat. I. DATEN



In Berlin demonstrierten tausende Demonstranten für mehr Datenschutz

Rund 7500 Demonstranten nahmen an dem Protestzug unter dem Motto "Freiheit statt Angst - Stoppt den Überwachungswahn" in Berlin teil.



5. DAAD Das Ende der Åra Bode



http://www.flickr.com/photos/topsy/204929063/

1. BUNDESLIGA Der BWB u







By Patrick Lauke

com



among the most imaginative architec-The National Aquatics Center lies on tural feats in recent memory. Critics have incessantly

described

a projects as bullish

the nation's budding - buildings are

Grand Central Term the great train halls Like Tempelhof, nal boasts à swee evokes the glame enclosing a surp

dragon. Yet i cedent Airport monument conceived Tim Griffith/PTW Architects Speer in ent ceremonial axis. t gateway Europe. Both are part o bile society that exte

lew China

although at times ten

fying in their aggress

scale, they also refl

the country's effor

give shape to an en

ing national identit Foster's airport t

nal, the world's li

is the purest exp.

of China's emb

the Modernist c

swooping form

suggests two

angs placed sid

has been com

http://www.flickr.com/photos/redux/409356158/





By Characterization of accounterplant Bring bacar on the day of web authors for create accessible content. However, with some and of the web accessible, cleans and use what accessibility is, why this important ba-whet accessibility is, why this important of web authors/developers to the buildity Authors buildity is and back indeveloper to the the project and how to choose up by An apportant of the second provides indeveloper to the buildity Authors indeveloper to the buildity Authors Rights Commission (RC). Web authors and developers need to be comes to developing their condexes when the comes to developing their condexes and evaluation. How ever, authors still need an actual inderstanding of the rationale undergraining understanding of the rationale undergraining and an accession in the second and undergraining inderstanding of the rationale undergraining and and an actual undergraining and accession in the second inderstanding of the rationale undergraining and and an accession in the second accession in the second inderstanding of the rationale undergraining and and accession in the second accession in the second inderstanding of the rationale undergraining and accession in the second accent basing and accession in the accession in the second accent basing and accession in the second accent basing and accession in the second accent basing and accent basing a



understanding of the rationale underpinning heckpoints and current best







understanding of the rationale underpinning

Finally, there

agents are a

particular

reiheit

Rund 7500 Demonstranten nahmen an dem Prot "Freiheit statt Angst – Stoppt den Überwachungs

and how

themselve



statt



ton straitet über k CTAG S/11 Obarna





EU AUF ZEIT ONLIN reiheit

Rund 7500 Demonstranten nahmen an dem Protestrag unter dem Mo "Freiheit statt Angst – Stoppt den Überwachungswahn" in Berlin teil.

B2-Particle Latation of the conservation of the part of the work according to the day of web automs the set of the set of

Finally, there has to be les. They need nts are available to th nd how to configure Telusi



ZEIT CONLINE DATENSCHUTZ

Taucondo domonetrioron für Bürgerrechte im Netz

Internet Date

Rund 7500 Demonstranten nahmen an dem Prote "Freiheit statt Angst – Stoppt den Überwachungsv

CRUZENSON UTRAMERI ÖRZÜNCI 2008 2008000000



ORF-Portal Futurezone kaufen aben als Kind gelennt, Teilen is PAN ALE EXCEPTOR "En End unequal ich halte en i IF & IPAD APPS Mere Freiheit im Ann Store

0

(in Datassuche Immobilien Autometici Joka Beineurgeboie

DATUM 11.9.2010 - 17:20 Uhv
1 - QUELLE ZEIT ONLINE, AFP, dos
q KOMMENTARE 4
 EMPPERLEN E-Mail verschicken | F

RTINEL ORUCKEN Dru

experteerde

Description of the end understanding of the rationale underpinning heckpoints and current best



a . to be Finally, there ha stude





From data to information.

Collect data and define your learning problem.

• Data preparation.

• Training a prediction model.

• Checking the performance of your model.

ZEIT CONLINE DATENSCHUTZ

Partnersuche Immobilien Automarkt Jobs Reiseangebote

STARTSEITE POLITIK WIRTSCHAFT MEINUNG GESELLSCHAFT KULTUR WISSEN DIGITAL STUDIUM KARRIERE LEBENSART REISEN AUTO

Internet Datenschutz Mobil Games

Anmelden | Registrieren

Suchen

DATENSCHUTZ

SPORT

Tausende demonstrieren für **Bürgerrechte im Netz**

Für einen besseren Arbeitnehmerdatenschutz und gegen die Gesundheitskarte: 130 Organisationen hatten zur Demonstration aufgerufen. Sie fürchten den Überwachungsstaat.



In Berlin demonstrierten tausende Demonstranten für mehr Datenschutz

Rund 7500 Demonstranten nahmen an dem Protestzug unter dem Motto "Freiheit statt Angst – Stoppt den Überwachungswahn" in Berlin teil.

netsenten vendten eich unter anderem zoren die Vellezsählung Dia Da

DATUM 11.9.2010 - 17:20 Uhr

- I- QUELLE ZEIT ONLINE, AFP, dpa
- C KOMMENTARE 4
- * EMPFEHLEN E-Mail verschicken | Facebook, Twitter, Buzz .
- ARTIKEL DRUCKEN Druckversion | PDF SCHLAGWORTE Datenschutz | Demonstration |
- Datensicherheit | Medienpolitik

NEU IM RESSORT

- I. DATENSCHUTZ Tausende demonstrieren für Bürgerrechte im Netz 2. FUTUREZONE Kurier darf ORF-Portal Futurezone kaufen
- 3. URHEBERRECHTE "Wir haben als Kind gelernt, Teilen ist
- out" 1. SPAM AUF FACEBOOK "Ein iPad umsonst, ich halte es in den Händen
- 5. IPHONE & IPAD APPS Mehr Freiheit im App-Store

NEU AUF ZEIT ONLINE

- I. NACHRUF Die Freie Bärbel Bohley ist tot 2. CDU Union streitet über konservatives Profil
- 3. GEDENKTAG 9/11 Obama warnt vor religiösen Ressentiments
- 1. BUNDESLIGA Der BVB und die Freiburger siegen 5. DAAD Das Ende der Åra Bode





ZEIT CONLINE | DATENSCHUTZ

STARTSEITE POLITIK WIRTSCHAFT MEINUNG GESELLSCHAFT KULTUR WISSEN DIGITAL STUDIUM KARRIERE LEBENSART REISEN AUTO

Internet Datenschutz Mobil Games

DATENSCHUTZ

SPORT

Tausende demonstrieren für Bürgerrechte im Netz

Für einen besseren Arbeitnehmerdatenschutz und gegen die Gesundheitskarte: 130 Organisationen hatten zur Demonstration aufgerufen. Sie fürchten den Überwachungsstaat.



Rund 7500 Demonstranten nahmen an dem Protestzug unter dem Motto "Freiheit statt Angst – Stoppt den Überwachungswahn" in Berlin teil.

Die Domonstranten wandten eich unter anderem zeren die Vellersählung

Anmelden | Registrieren

Partnersuche Immobilien Automarkt Jobs Reiseangebote

- KOMMENTARE 4
 EMPFEHLEN E-Mail verschicken | Facebook, Twitter,
 Buzz
- ARTIKEL DRUCKEN Druckversion | PDF
 SCHLAGWORTE Datenschutz | Demonstration |
 Datensicherheit | Medienpolitik

NEU IM RESSORT I. DATENSCHUTZ Tausende demonstrieren für Bürgerrechte

im Netz 2. FUTUREZONE Kurier darf ORF-Portal Futurezone kaufen 3. URHEBERRECHTE "Wir haben als Kind gelernt, Teilen ist

gut" 5. SPAM AUF FACEBOOK "Ein iPad umsonst, ich halte es in den Händen" 5. IPHORE & IPAD APPS Mehr Freiheit im Apo-Store

NEU AUF ZEIT ONLINE I. NACHRUF Die Freie - Bärbel Bohley ist tot 2. CDU Union streitet über konservatives Profil

 S. GEDENKTAG 9/11 Obama warnt vor religiösen Ressentiments
 BUNDESLIGA Der BVB und die Freiburger siegen
 S. DAAD Das Ende der Ära Bode

ANZEIGE

• Remove noise.

ZEIT CONLINE | DATENSCHUTZ

STARTSEITE POLITIK WIRTSCHAFT MEINUNG GESELLSCHAFT KULTUR WISSEN DIGITAL STUDIUM KARRIERE LEBENSART REISEN AUTO

Internet Datenschutz Mobil Games

DATENSCHUTZ

Tausende demonstrieren für Bürgerrechte im Netz

Für einen besseren Arbeitnehmerdatenschutz und gegen die Gesundheitskarte: 130 Organisationen hatten zur Demonstration aufgerufen. Sie fürchten den Überwachungsstaat.



Rund 7500 Demonstranten nahmen an dem Protestzug unter dem Motto "Freiheit statt Angst – Stoppt den Überwachungswahn" in Berlin teil.

Die Demonstranten wandten eich unter anderem zonen die Vellezeihlung

Anmelden | Registrieren

Partnersuche Immobilien Automarkt Jobs Reiseangebote

- * EMPFEHLEN E-Mail verschicken | Facebook, Twitter, Buzz ...
- ARTIKEL DRUCKEN Druckversion | PDF
 SCHLAGWORTE Datenschutz | Demonstration |
 Datensicherheit | Medienpolitik

NEU IM RESSORT I. DATENSCHUTZ Tausende demonstrieren für Bürgerrechte

im Netz 2. FUTUREZONE Kurier darf ORF-Portal Futurezone kaufen 3. URHEBERRECHTE "Wir haben als Kind gelernt, Teilen ist

gut" **§ SPAM AUF FACEBOOK** "Ein iPad umsonst, ich halte es in den Händen"

5. IPHONE & IPAD APPS Mehr Freiheit im App-Store

NEU AUF ZEIT ONLINE

5. DAAD Das Ende der Ära Bode

NACHRUF Die Freie - Bärbel Bohley ist tot
 COU Union streifet über konservatives Profil
 GEDENKTAG 9/11 Obama wamt vor religiösen
 Ressentiments
 BUNDESLIGA Der BVB und die Freiburger siegen

avzeige

• Remove noise.

• Convert text to vectors.

From texts to vectors

If we looked at two words only:









Binary bag of words

- Imagine a n-dimensional space.
- Each dimension = one possible word in texts.
- Entry in vector is one, if word occurs in text.

$$b_{i,j} = \begin{cases} 1 \forall x_i \in d_j \\ 0 \, else \end{cases}$$

- Problem:
 - Number of word occurrences not accounted for.

Term Frequency

- Imagine a n-dimensional space.
- Each dimension = one possible word in texts.
- Entry in vector equal to the words frequency.

$$b_{i,j} = n_{i,j}$$

- Problem:
 - Common words dominate vectors.

TF with stop wording

- Imagine a n-dimensional space.
- Each dimension = one possible word in texts.
- Filter stopwords.
- Entry in vector equal to the words frequency.

$$b_{i,j} = n_{i,j}$$

- Problem:
 - Common and uncommon words with same weight.

TF- IDF

- Imagine a n-dimensional space.
- Each dimension = one possible word in texts.
- Filter stopwords.
- Entry in vector equal to the weighted frequency.

$$b_{i,j} = n_{i,j} \times \log\left(\frac{|D|}{|[d:t_i \in d]|}\right)$$

- Problem:
 - Long texts get larger values.

Normalized TF- IDF

- Imagine a n-dimensional space.
- Each dimension = one possible word in texts.
- Filter stopwords.
- Entry in vector equal to the weighted frequency.
- Normalize vectors.

$$b_{i,j} = \frac{n_{i,j}}{\sum_{k} n_{k,j}} \times \log\left(\frac{|D|}{|[d:t_i \in d]]}\right)$$

- Problem:
 - Additional domain knowledge ignored.

Reality

- There are a few more words in news.
- Use all relevant features/ signals available.
 - Words.
 - Header fields.
 - Characteristics of publishing url.
 - •
- Usually pipeline of feature extractors.

From data to information.

Collect data and define your learning problem.

• Training a prediction model.

• Checking the performance of your model.

Step 2: Similarity









Step 3: Clustering












Reality

• Seed selection.

• Choice of initial k.

• Continuous updates.

• Regular addition of clusters.

Discover groups of similar items

Canopy.
 Dirichlet based.

k-Means.
 Spectral clustering.

• Fuzzy k-Means. • Others upcoming.

From data to information.

Collect data and define your learning problem. Data preparation. Training a prediction model.

• Checking the performance of your model.

Evaluation

• Compare against gold standard.

• Use quality measures.

• Manual inspection.

From data to information.

Collect data and define your learning problem. Data preparation. Training a prediction model. Checking the performance of your model.

http://www.flickr.com/photos/generated/943078008/

ttp://www.flickr.com/photos/eschipul/4160817135/



What else does Mahout have to offer.

Identify dominant topics

• Given a dataset of texts, identify main topics.

Algorithms: Parallel LDA

- Examples:
 - Dominant topics in set of mails.
 - Identify news message categories.

Assign items to defined categories.

• Given pre-defined categories, assign items to it.



ickr.com/photos/63056612@N00/155554663/ By freezelight, http://

SPAM

SPAM

SPAN

SPAM

SPAM

SPAN

SPAN SPAN

SPAM

Hormel

SPAI

SPAN

SPA

SPA

SPAM

SPAM

SPA

Leving .

Horme

SPAM

SPAM

SPAR

Google images oakland

SafeSearch: Moderate V

Related searches: oakland raiders

Images - Hide options

> Any size Medium Large Icon Larger than... Exactly

> Any type Face Photo Clip art Line drawing

> Any color Full color Black and white Specific color

Oakland Airport 625 x 471 - 103k - jpg visitingdc.com



Search images

Oakland Ranks Fifth 538 x 359 - 44k - jpg bayareahomegirl.com



Advanced Image Search

Oakland 900 x 600 - 171k - jpg globalsecurity.org



Result

OAKLAND First 400 x 400 - 27k bayassociation.org



Oakland Gaudy Lexus 2 500 x 340 - 73k - jpg lexusenthusiast.com



The Oakland Baseball 450 x 305 - 39k - jpg sportsbusinesssims.com



Oakland 600 x 320 - 48k - jpg webpages.scu.edu



Learn more about 550 x 366 - 56k - jpg tripadvisor.com

Google images oakland

SafeSearch: Moderate V

Related searches: oakland raiders

Images > Face - Hide options

> Any size

Medium Large

Icon Larger than ...

Exactly

Any t Face Photo

Clip an Line drawing

> Any color Full color Black and white Specific color

Reset options

RAIDERS

Jakland, CA 94621 513 x 545 - 13k - gif nflfootballstadiums.com



Search images

All Graphics » 262 x 278 - 43k - gif coolchaser.com



Advanced Image Search

Oakland Sideshow and 720 x 480 - 44k channels.com



Detroit Tigers v 594 x 396 - 51k zimbio.com



oakland@coe.ufl.edu 379 x 471 - 110k - jpg coe.ufl.edu



Results 1 - 20 of about 1,400,000 (0.

Oakland. by Davey D 359 x 512 - 26k - jpg sfbayview.com



#20 of the Oakland 467 x 594 - 86k zimbio.com



Detroit Tigers v 443 x 594 - 74k zimbio.com



Dallas Cowboys v 594 x 404 - 60k zimbio.com



Distributed by Tubemogul. The 720 x 480 - 18k channels.com



Assign items to defined categories.

• Naïve Bayes. • Random forests.

- Complementary naïve
 Logistic regression.
 - HMM for sequences.

Recommendation mining.

• Collaborative filtering.



Show most relevant ads

ALUMINIUM Baseballschläger 30' American Baseball

von Outdoor 4 You - Shop

· 神宮古古古 문 (4 Kundenrezensionen) Mehr zu diesem Artikel

Preis: EUR 17,58

N DOLLAR DE LA COLOR

Auf Lager. Verkauf und Versand durch NORMANI.

Noch 5 Stück auf Lager. 4 neu ab EUR 17,58



Show most relevant ads

erwandte Su	chbegriffe: <mark>clearasil pickelstift, world of warcraft.</mark>
este Ergebniss	e Zurück Seit
10	Clearasil 44161 Tiefenreinigung Antibakterielle Reinigungspads, 60e
Cibar	Lieferung bis Samstag, 22. März: Bestellen Sie innerhalb der nächsten 23 Stunden Kostenlose Lieferung möglich:
	大学会会会 (1) Drogerie & Bad: Alle 13 Artikel ansehen
2	Clearasil Ultra Anti-Pickel Reinigungspads, 65 Stück von Clearasil (Ba
	Neu kaufen: EUR 7,99
Dearan	Gewöhnlich versandfertig in 1 bis 3 Wochen. Kostenlose Lieferung möglich.
	Drogerie & Bad: Alle 13 Artikel ansehen
NARER	World of WarCraft: Wrath of the Lich King (Add-on) von Vivendi Unive Vista / XP)
C.	Neu kaufen: EUR 39,99

Vorbestellbar Kostenlose Lieferung möglich. Games: Alle Artikel ansehen

Show most relevant ads



~ <u>Otis Gospodnetic</u> (Author), <u>Erik Hatcher</u> (Author)

List Price: \$44.95

Price: \$29.67 & this item ships for FREE with Super Saver Shipping. Details You Save: \$15.28 (34%)

In Stock.

Ships from and sold by Amazon.com. Gift-wrap available.

23 new from \$19.99 15 used from \$17.11



Frequently Bought Together

Customers buy this book with Building Search Applications: Lucene, LingPipe, and Gate by Manu Konchady



Customers Who Bought This Item Also Bought





Hadoop: The Definitive Guide by Tom White



Hibernate Search in Action by Emmanuel







Collective Intelligence in Action by Satnam Alag

Recommending places

http://www.flickr.com/photos/sebastian_bergmann/1244514498

http://www.flickr.com/photos/jfclere/4061801735













Thanks to Falko Menge for the pictures of Brussels.



Recommending people











Recommendation mining.

- Online collaborative filtering on single machine.
- Offline Map/Reduce based version.
- Content similarity can be integrated.

• Based on former Taste project.

Frequent pattern mining

 Given groups of items, find commonly cooccurring items.

- Examples:
 - In shopping carts find items bought together.
 - In query logs find queries issued in one session.





rypto/3201254932/sizes/l/

By libraryman, http://www.flickr.com/photos/libraryman/78337046/sizes/l/

By quinnanya, http://www.flickr.com/photos/quinnanya/2806883231/



By crypto, http://www.flickr.com/photos/crypto/3201254932/sizes/l/

By libraryman, http://www.flickr.com/photos/libraryman/78337046/sizes/l/

Requirements to get started

March 14, 2009 by Artful Magpie http://www.flickr.com/photos/kmtucker/3355551036/





•• AWS	··· Products	· Developers	Community	Y Support	× Account
Products & Services	Amazon	Elastic Compu	te Cloud (Amaz	on EC2)	
Amazon EC2 Details	Amazon Elastic resizable comp	c Compute Cloud (Amazon ute capacity in the cloud.	EC2) is a web service that It is designed to make web-	provides Sign Up F	or Amazon EC2 🕥
EC2 Overview	computing eas	ier for developers.			
FAQs	Amazon EC2's	simple web service interfa	ice allows you to obtain and	Č.	
Amazon EC2 SLA configure capacity with minimal friction. It provides you with complete control of your computing resources and late you run on Amazon's preven					
EC2 Instance Types	computing env	obtain			

Amazon Elastic MapReduce

Amazon Elastic MapReduce is a web service that enables businesses, researchers, data analysts, and developers to easily and cost-effectively process vast amounts of data. It utilizes a hosted Hadoop framework running on the web-scale infrastructure of Amazon Elastic Compute Cloud (Amazon EC2) and Amazon Simple Storage Service (Amazon S3).

Using Amazon Elastic MapReduce, you can instantly provision as much or as little capacity as you like to perform data-intensive tasks for applications such as web indexing, data mining, log file analysis, machine learning, financial



(Thanks to Thilo for helping set up the cluster, Thanks to packet and masq for two of the three machines.)





Why go for Apache Mahout?

Jumpstart your project with proven code.

January 8, 2008 by dreizehn28 http://www.flickr.com/photos/1328/2176949559

Discuss ideas and problems online.

November 16, 2005 [phil h] http://www.flickr.com/photos/hi-phi/64055296









Become a committer.















Sebastian Schelter Jake Mannix Benson Margulies Robin Anil David Hall AbdelHakim Deneche Karl Wettin Sean Owen Grant Ingersoll Otis Gospodnetic Drew Farris Jeff Eastman Ted Dunning Isabel Drost



Become a committer: Of Apache Mahout







Emeritus:

Niranjan Balasubramanian Erik Hatcher Ozgur Yilmazel Dawid Weiss
-user@.apache.org *-dev@*.apache.org



Image by: Patrick McEvoy

Interest in solving hard problems. Being part of lively community. Engineering best practices.

Bug reports, patches, features. Documentation, code, examples.

Thanks to Tim Lossen et. al for taking amazing pictures of the conf.

ho

10

Herzlich Willkommen!

l love lelvetik a

Berlin Buzzwords 2011

Search/ Store/ Scale

May/ June 2011

Thanks to Tim Lossen et. al for taking amazing pictures of the conf.

*-user@hadoop.apache.org *-dev@hadoop.apache.org



Image by: Patrick McEvoy

Interest in solving hard problems. Being part of lively community. Engineering best practices.

Bug reports, patches, features. Documentation, code, examples.