

Course Introduction and Overview

Internet Analytics (COM-308)

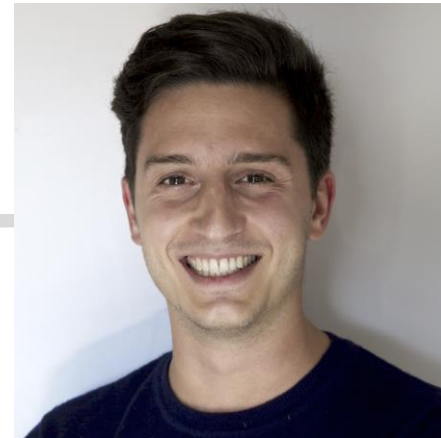
Prof. Matthias Grossglauser
School of Computer and Communication
Sciences



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Team

- Instructor:
 - Matthias Grossglauser
- Assistants:
 - Victor Kristof ●
 - William Trouleau ●
- Team of student-assistants:
 - Ahmad: homework sessions
 - Olivier, Orcun, Alexis: labs



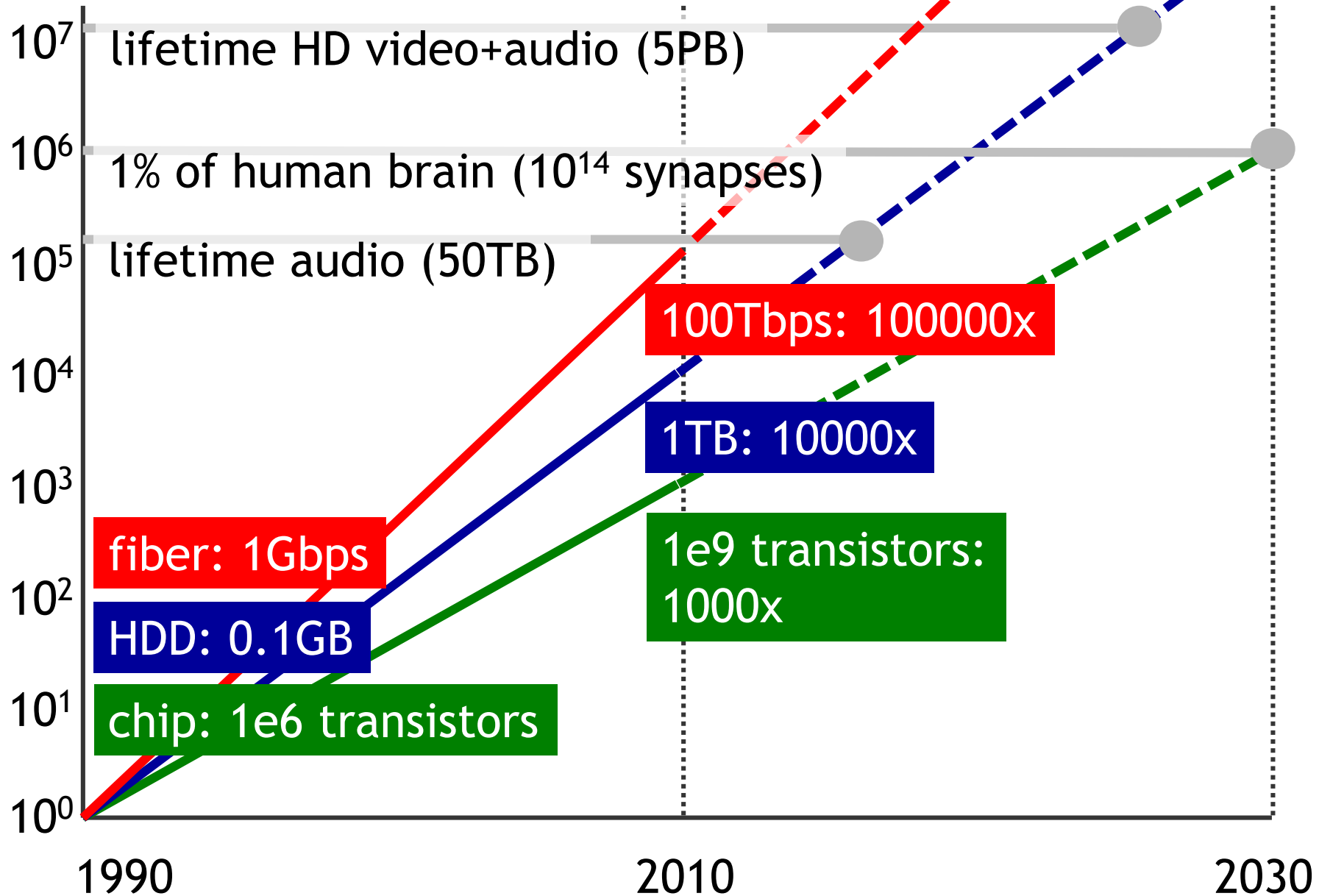
Logistics

- URL:
 - Moodle module “Internet Analytics”: slides, labs, Q&A
- Labs:
 - BC07-08
 - Wed morning 8:15-10:00
 - Groups of 3 - self-organize & register in moodle (DL: end of next week)
 - Spark account setup: please register for the class by **Tuesday 17:00h**
- Homeworks
 - Wednesday after lab (BC01)
- Midterm exam: **TBD by next week**
- Final exam: **Wed May 30, 08:15-11:00**

Grading

- Midterm exam: 30%
- Final exam: 50%
- Homeworks: not graded
 - But strongly recommended to do them regularly
- Labs: 20%
 - Lab-0 not graded
 - Hand-ins (code, plots, interpretations,...)
 - Deadline for lab n (usually) at the start of lab $n + 1$
 - Mini-interviews to check your understanding

Technology evolution: exponential trends



Limited prediction horizons

"This 'telephone' has too many shortcomings to be seriously considered as a means of communication." – Western Union

"Computers in the future may have only 1,000 vacuum tubes and perhaps only weigh 1 1/2 tons." – Popular Mechanics

"There is no reason for any individual to have a computer in their home." – Ken Olson, President DEC

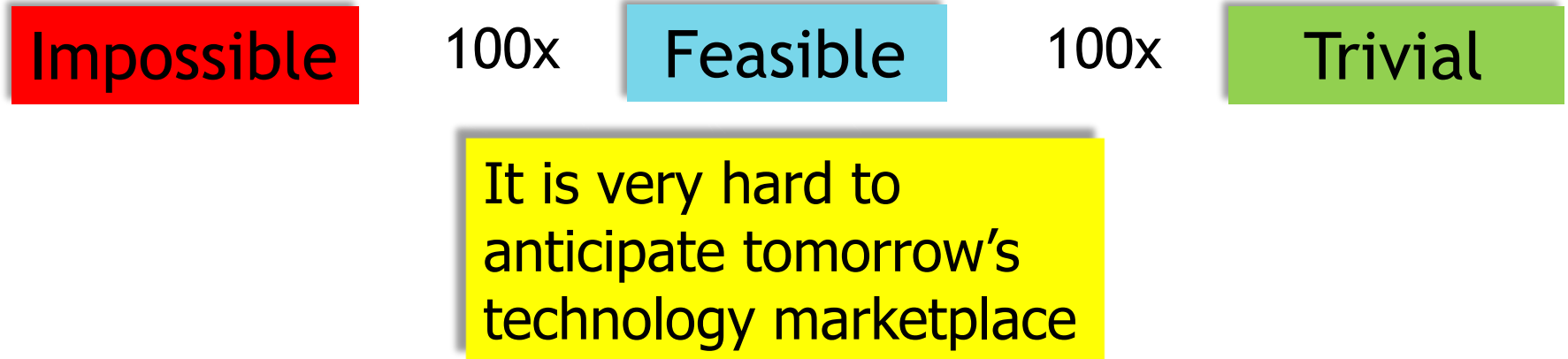
1850

1900

1950

2012

Exponential technology evolution



- Technology fundamentals
 - Several surprising exponential scaling laws
 - Stable and predictable over decades
 - What applications in 10-20 years: no idea
 - Clear trend: measuring, storing, analyzing everything!

Evolution of interfaces & user data

Social networks & processes, sharing patterns, preferences,...

Fine-grained mobility and behavior patterns

Coarse-grained mobility

Organization of static information, access patterns

smartphones

social / web 2.0

mobile

web

1990

2000

2015

What is this class about?

- Web and mobile services
 - Explosive growth: social networking and social media; messaging; blogs; location & navigation; ...
- Pre-cloud engineering challenges:
 - Networking: connectivity, routing, traffic engineering, intrusion detection,...
 - Data-center design: databases, server farms, traffic monitoring, energy, hosting,...
- Post-cloud engineering challenges:
 - Cloud: outsourcing of many “lower-layer functions”: connectivity, storage, computing, security,...
 - Data science: drives user functionality & user experience; monetization (ads, etc.); management (capacity, etc.),...

What is this class about?

- Social web, online social networks, mobile:
 - Explosion of user data
 - Increasing correlation of user data: more than sum of its parts
 - Demand for “information engineers” or “data scientists”
- Data: a huge variety
 - What are the main user data types?
 - What underlying models describe them?
- Function/application-oriented
 - How is data turned into decisions, actions, insights?
- Working with real data

What is this class about?

- Real data:
 - Degree of realism
 - Interesting in its own right
 - Real-world challenges: gaps, errors, lack of documentation,...

What this class is not about?

- Little on business models, policy & legal issues
- Little coverage of systems issues
 - Cloud architectures, energy,...
- Not an exhaustive ML class
 - E.g., no self-contained treatment of “classical” topics like regression, SVM, deep learning - but introduced as needed
- Criterion in choice of topics: real-world relevance of...
 - Data & models
 - Functions & applications

Matrix of data and functions

	Networks	Ratings	Document	Corpus	Streams
Characterize Model	Small worlds, scale-free			Topic models	Counts, moments
Predict Infer	Link prediction				
Rank	PageRank	Collaborative filtering	Content-based recommend		
Filter			Spam filtering		
Search Retrieve				Keyword search, similar docs	
Associate Summarize	Community detection			Clustering	Random projections

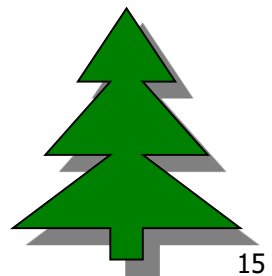
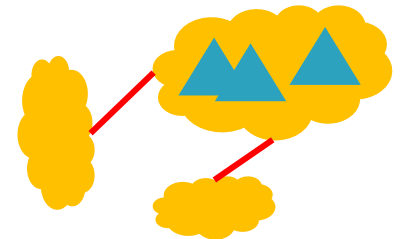
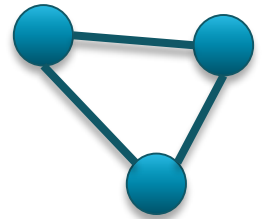
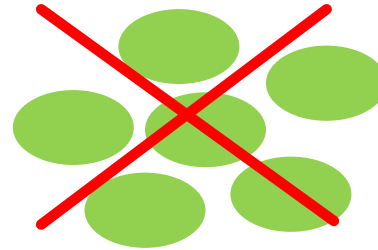
Overview: social and info networks

- Social and information networks
 - How are people connected?
 - How is information connected?



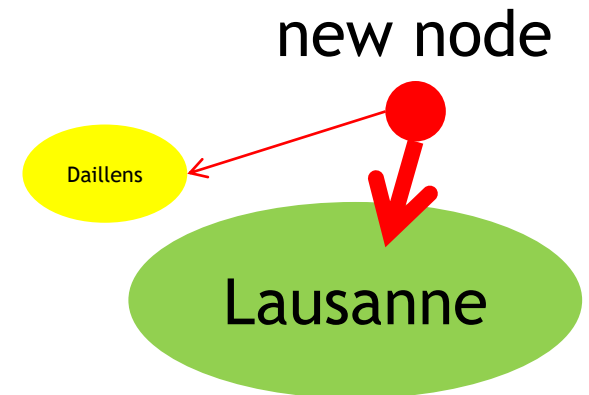
Networks: structure

- Structure:
 - Single snapshot in time
 - Generic properties?
- 1: Giant Component
 - Almost all pairs of nodes are reachable
- 2: Clustering
 - Many triangles
- 3: Strong and weak ties
 - Interconnected sub-communities
 - Your acquaintances are more important than your friends (for finding a job ;-)
- 4: Short distances
 - Six degrees of separation



Networks: evolution

- Evolution:
 - How does the network change over time?
- Herding behavior:
 - We tend to copy behavior of those around us
 - Benefitting from each others' decisions
- 1: Information cascades
 - Sequences of wrong decisions
- 2: Preferential attachment
 - Skewed degree distribution
 - The rich get richer
- 3: Friendship paradox
 - Your friends have more friends than you

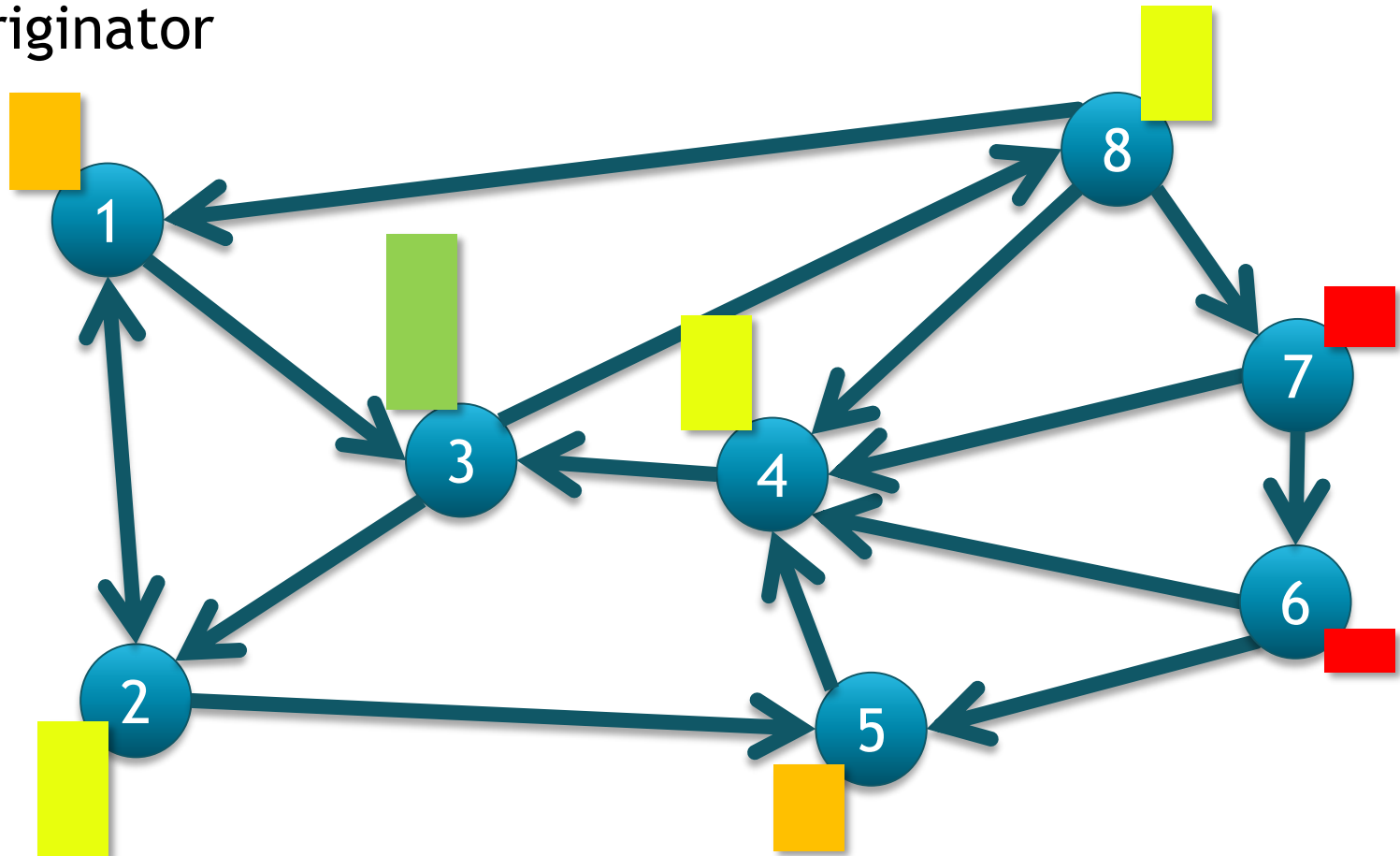


Networks: processes

- Processes:
 - Nodes have state and influence each other
 - How does net structure influence processes?
- 1: Epidemics
 - How does a rumor spread through a social network?
 - How does a disease progress, and who should be vaccinated?
- 2: Sampling
 - Very large network: how to estimate properties without visiting all the nodes?

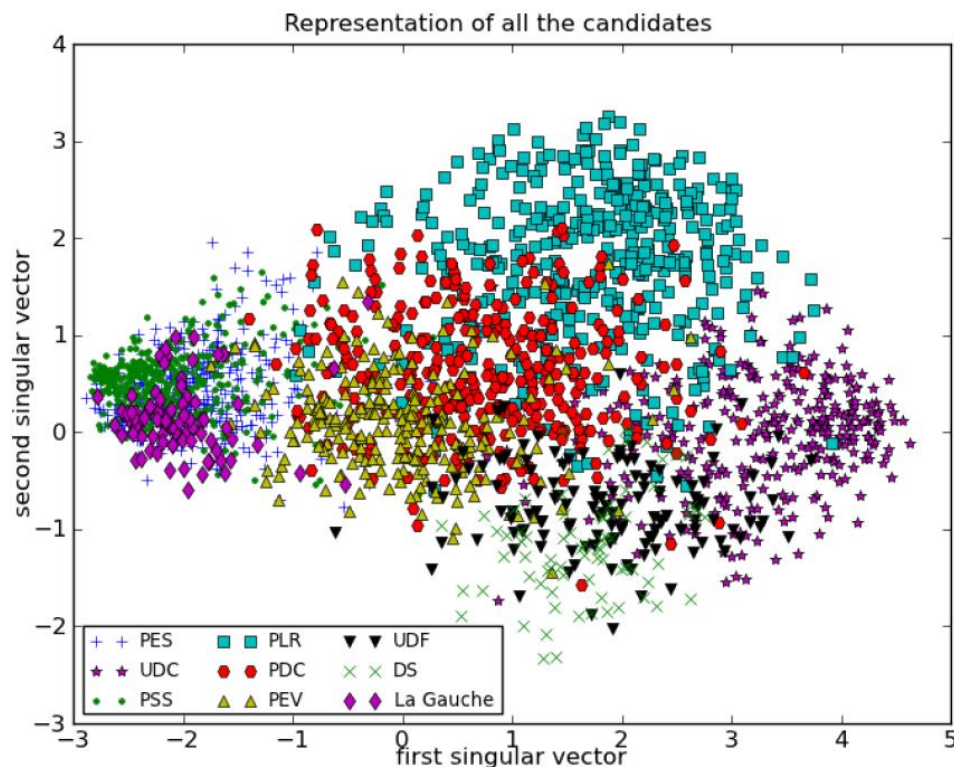
Overview: relevance and filtering

- PageRank:
 - A hyperlink “endorses” the target
 - An endorsement depends on the “relevance” of the originator



Overview: dimensionality reduction

- Raw data:
 - Often high-dimensional
 - But has “structure” = low-dimensional signal + noise
- Challenge:
 - How to find low-dimensional structure?
- Applications:
 - Visualizing
 - Explaining
 - Modeling
- Example:
 - SmartVote dataset on political candidates



Overview: recommender systems

[matthias's Amazon.com](#) | [Today's Deals](#) | [Gift Cards](#) | [Help](#)

Shop by Department ▾

Search All ▾

[Your Amazon.com](#) | [Your Browsing History](#) | [Recommended For You](#) | [Amazon Betterizer](#) | [Improve Your Recommendations](#) | [Your Profile](#) | [Learn More](#)

[Your Amazon.com](#) > **Recommended for You**
(If you're not matthias grossklauser, click here.)

Just For Today
[Browse Recommended](#)

Recommendations
[Amazon Instant Video](#)
[Amazon MP3 Store](#)
[Appliances](#)
[Appstore for Android](#)
[Arts, Crafts & Sewing](#)
[Automotive](#)
[Baby](#)
[Beauty](#)
[Books](#)
[Books on Kindle](#)
[Camera & Photo](#)
[Cell Phones & Accessories](#)
[Clothing & Accessories](#)
[Computers](#)
[Electronics](#)
[Grocery & Gourmet Food](#)
[Health & Personal Care](#)
[Home & Kitchen](#)
[Home Improvement](#)

These recommendations are based on [items you own](#) and more.

view: [All](#) | [New Releases](#) | [Coming Soon](#)

1.



[Curb Your Enthusiasm: The Complete Eighth Season](#)
DVD ~ Larry David (June 5, 2012)
Average Customer Review: ★★★★★ ☒ (88)
In Stock

List Price: ~~\$39.98~~
Price: **\$25.77**
[16 used & new](#) from \$25.77

☐ I own it ☐ Not interested ☒ ★★★★★ Rate this item
Recommended because you purchased [Curb Your Enthusiasm: The Complete Sixth Season](#) and more ([Fix this](#))

2.



[Probabilistic Graphical Models: Principles and Techniques \(Adaptive C](#)
by Daphne Koller (July 31, 2009)
Average Customer Review: ★★★★★ ☒ (17)
In Stock

List Price: ~~\$99.00~~
Price: **\$95.04**
[81 used & new](#) from \$64.98

☐ I own it ☐ Not interested ☒ ★★★★★ Rate this item
Recommended because you added [Introduction to Information Retrieval](#) to your Shopping Cart and more ([Fix this](#))

Overview: recommender systems

- Content-based recommenders

item 1:
“Plane hijacked...”

item 2:
“soccer game...”

item 3:
“swiss skiers win...”

item 4:
“50.3% vote yes...”



News item 1

News item 2

News item 3

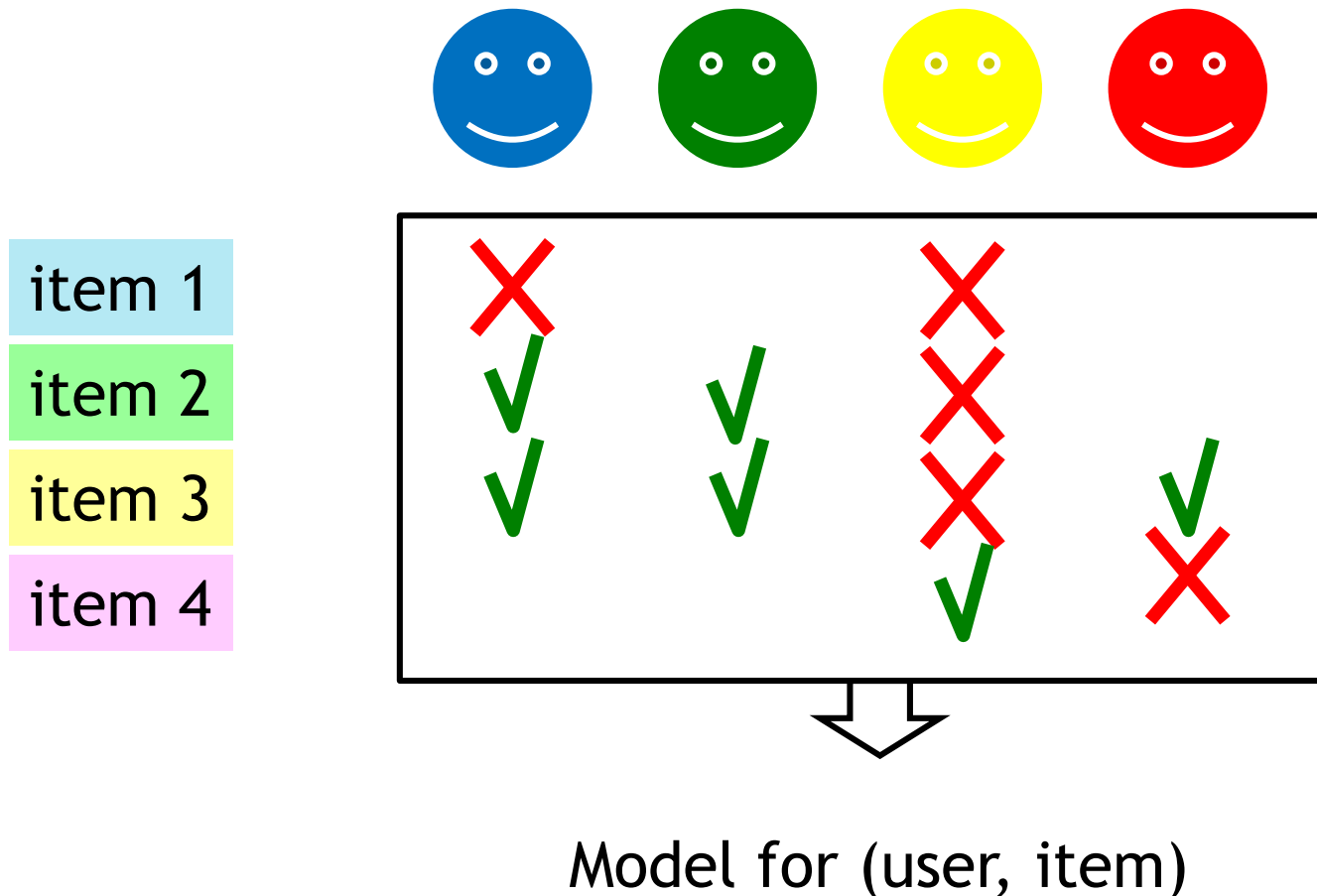
News item 4



Model for
(user, content)

Overview: recommender systems

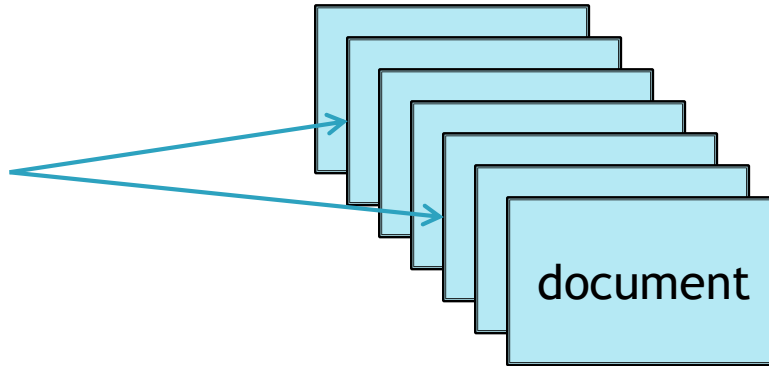
- Collaborative filtering-based recommenders



Overview: search and retrieval

- Given a query, how to find best matches?

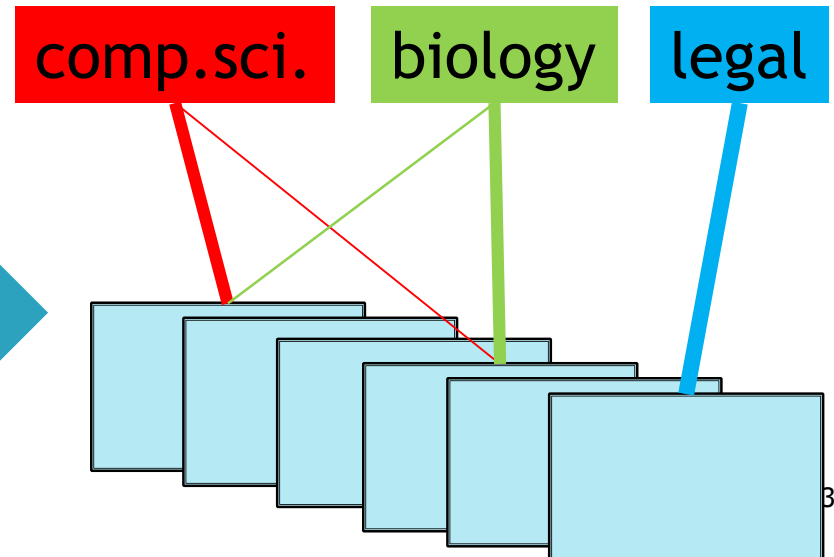
“internet analytics”



- Without a query, how to describe a corpus?

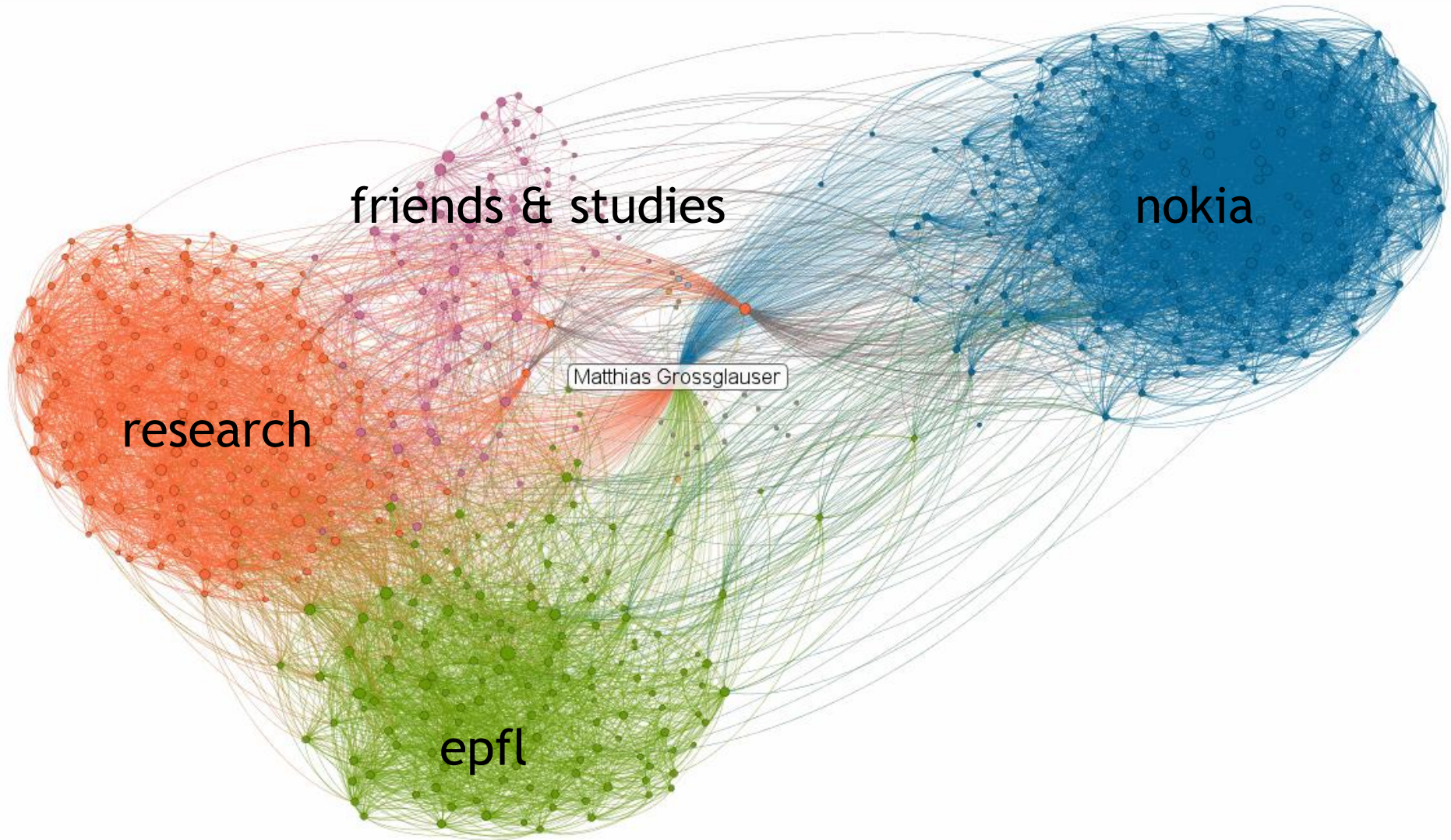


topic models



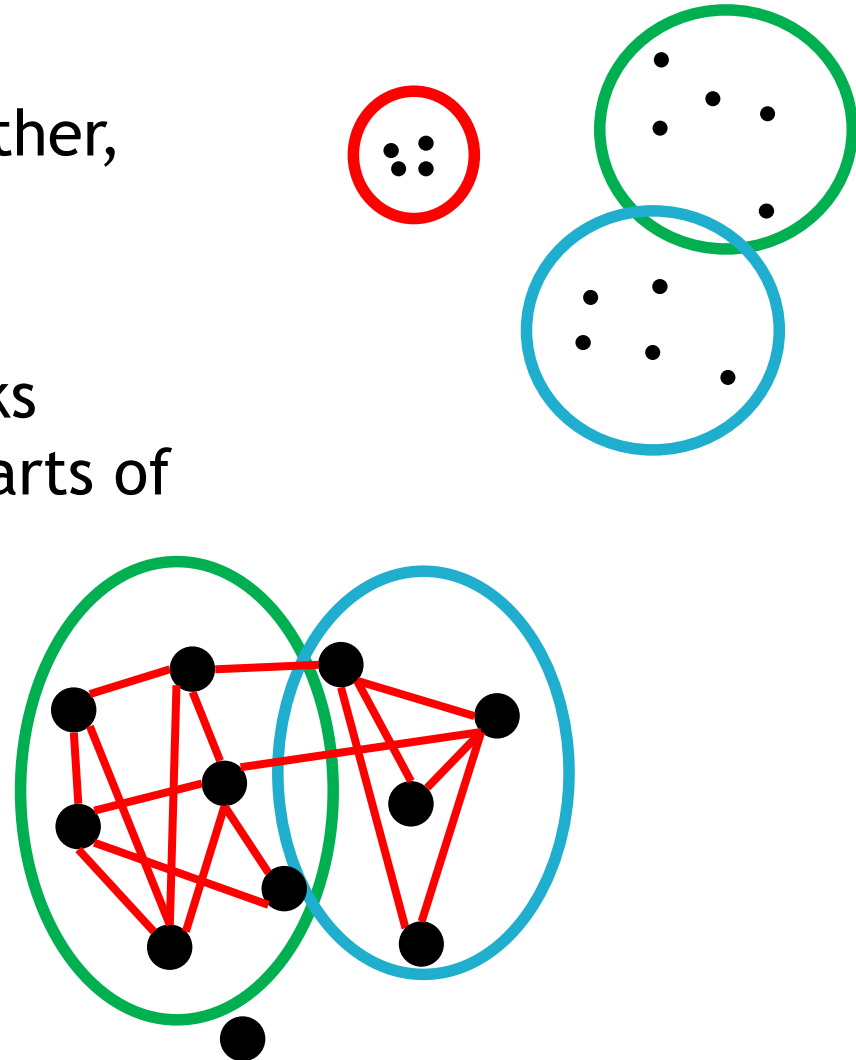
Overview: clusters, groups, communities

LinkedIn Maps Matthias Grossglauser's Professional Network
as of February 16, 2013



Overview: clusters, groups, communities

- Another type of structure
- Cluster:
 - Set of points close to each other, and far from other points
- Community:
 - A set of nodes with more links among them than to other parts of network



Overview: streams

- Internet backbone router
- Order of magnitude:
 - 100s of interfaces at 10s of Gbps
 - = several billion pkts/sec!
- Traffic analysis app to detect DDoS attack:
 - How many *different* (unique) source IP addresses in a minute?
 - If too large -> suspicious (fake addresses)!



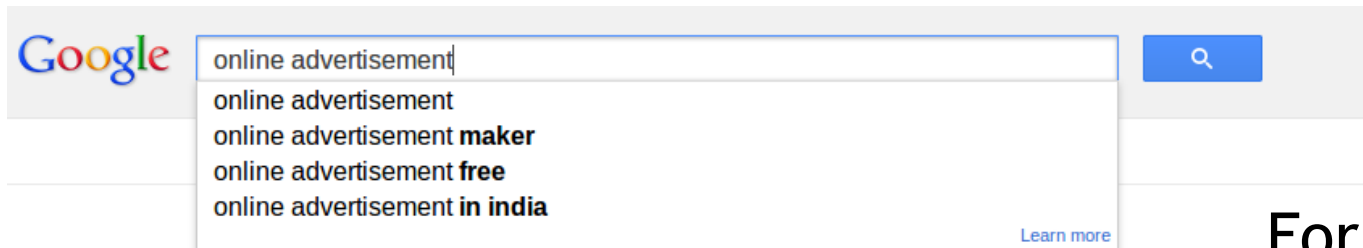
Overview: streams

- Computing statistics with sub-linear memory
- Example:
 - n numbers: how many unique values k ?



- How to solve with $\theta(n)$ memory?
 - Keep every value in some efficient data structure; compare & count
- How to solve with $o(n)$ memory?
 - Cannot solve exactly
- Streaming algorithms: ∞ data, finite memory
 - Approximation
 - (Pseudo-)randomization

Overview: ad auctions



Ads related to **online advertisement** ⓘ

Online Email Advertising - VerticalResponse.com
www.verticalresponse.com/FreeTrial
Trusted by 500k+ for Reliable Email Marketing Since 2001. Try it Free!
238 people +1'd or follow VerticalResponse
Flexible Pricing Options - Features - Pay As You Go - Start Free Trial Now

Online Advertising - Create your online marketing strat.
www.agencevirtuelle.com/OnlineMarketing
AdWords, SEO, Mobile, Social Media.

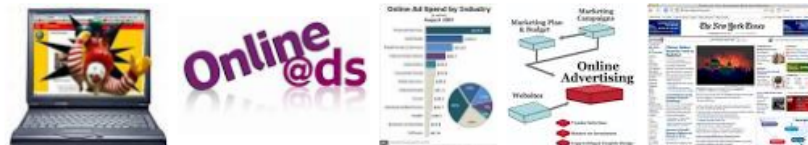
For each search, this table of “sponsored search results” is the result of an online auction

Online advertising - Wikipedia, the free encyclopedia
en.wikipedia.org/wiki/Online_advertising

Online advertising, also known as **online advertisement**, internet marketing, online marketing or e-marketing, is the marketing and promotion of products or ...
[History of online advertising - Competitive advantage over ... - Online advertisement](#)

Online Advertising: How to Do It Right | Small Business Trends
smallbiztrends.com/2010/11/online-advertising-how-to-do.html
Nov 4, 2010 – Helpful tips on using **online advertising** for small businesses.

Images for online advertisement - Report images



Internet advertising: The ultimate marketing machine | The Economist
www.economist.com/node/7138905

July 6, 2006 – IN TERMS of efficiency, if not size, the **advertising** industry is only now

Overview: ad auctions

- Online advertisement:
 - Key business model for the (consumer) internet & mobile
- Keyword auctions:
 - Each request to google -> auction
 - Advertiser bid for (keyword, bid)
- Generalized Second Price (GSP) auction

Class Schedule (1)

1	Introduction and Overview	<ul style="list-style-type: none">- class overview- logistics
2	Social and Information Networks 1: Structure	<ul style="list-style-type: none">- intro: social networks, web, social web- social networks, key properties, effects & metrics- giant component, relation to tree percolation- strong/weak ties, clustering- "everything is close": distances; random graphs
3	Social and Information Networks 2: Evolution	<ul style="list-style-type: none">- "the rich get richer": power laws, cumulative advantage, pref attachment- my friends have more friends than i-phenomenon- link prediction
4	Social and Information Networks 3: Processes	<ul style="list-style-type: none">- processes: epidemics, cascades- sampling a network: undirected vs directed- random walks on networks, mixing times, spectral properties

Class Schedule (2)

5	Ranking	<ul style="list-style-type: none">- Intro web structure- PageRank algorithm- Large-scale computation- HITS variant: hubs and authorities
6	Dimensionality Reduction	<ul style="list-style-type: none">- intro: "finding hidden structure", visualization- PCA and derivatives- singular-value decomposition- manifold unwrapping
7	Recommender Systems 1	<ul style="list-style-type: none">- Collaborative filtering- graph-based, item-item vs user-user- spectral/matrix completion- case study: netflix
8	Recommender Systems 2	<ul style="list-style-type: none">- Intro: "the long tail"- Applications, models- TF-IDF- kNN classifier- Naïve Bayes classifier

Class Schedule (3)

9	Clusters and Communities	<ul style="list-style-type: none">- Gaussian Mixture Model (GMM)- EM algorithm- Communities and modularity measures- Louvains clustering algorithm
10	Text Search and Retrieval 1	<ul style="list-style-type: none">- Latent Semantic Indexing (LSI)- Intro: Bayesian networks- Gibbs sampling
11	Text Search and Retrieval 2	<ul style="list-style-type: none">- Probabilistic LSI (pLSI)- Latent Dirichlet Allocation (LDA)
12	Streams	<ul style="list-style-type: none">- Intro- motivating apps- streaming model
13	Internet Ad Economy & Online Auctions	<ul style="list-style-type: none">- intro: sponsored search, keyword auctions- VCG and Generalized Second Price (GSP)- incentive compatibility

Caveats (1)



- Breadth
 - Intersection of data mining, machine learning, network science, statistics, algorithms
- No textbook
 - Combination of several textbooks & other sources
 - The slides + your notes are the course support
- Focus on real applications & data
 - Details often unknown: trade secrets
 - E.g., google practice vs original PageRank
- Lecture/lab overlap
 - Lectures+homeworks: models, theory, background
 - Labs: try it out on real data
 - Overlap is deliberately partial: you learn something new in labs

Caveats (2)



Previously	2017, 2018
Hadoop	Spark
Labs focused on map-reduce model	Lab focus more general, map-reduce “under the hood”
Custom cluster built and maintained in my lab	Uses IC cloud infrastructure
Labs in Java	Labs in Python
~30 students	~60 students

- Implications:
 - Labs had to be redesigned from scratch
 - From tightly scripted labs to more open-ended labs
 - We cannot help you (too much) with low-level Python programming issues
 - Fewer but bigger labs
 - Groups of 3 rather than 2

Summary

- Problems & data from the real world...
- But enough theory, models to understand the foundations
- Required background:
 - Python & eclipse
 - Linear algebra
 - Probability & statistics
 - Algorithms
- Next:
 - Wed 08:15h: Tutorial on Spark and using infrastructure
 - No HW session this week yet
- Reminder: sign up for class in ISA by tomorrow!