



# Introduction

# SOCIAL

# MEDIA

# MINING



## **Dear instructors/users of these slides:**

Please feel free to include these slides in your own material, or modify them as you see fit. If you decide to incorporate these slides into your presentations, please include the following note:

R. Zafarani, M. A. Abbasi, and H. Liu, *Social Media Mining: An Introduction*, Cambridge University Press, 2014.  
Free book and slides at **<http://socialmediamining.info/>**

or include a link to the website:  
**<http://socialmediamining.info/>**

# Facebook

The image shows a screenshot of Mark Zuckerberg's Facebook profile. The header includes the Facebook logo, a search bar, and navigation links for Home, Profile, and Account. The profile picture is a large photo of Mark Zuckerberg. To the left of the profile picture is a sidebar with links to Wall, Info, Photos (826), Questions, and Family. The main content area displays Mark's name, a bio stating he has worked at Facebook, studied Computer Science at Harvard University, lives in Palo Alto, California, and was born on May 14, 1984. Below the bio are several small photos. The 'Education and Work' section lists his employers: Facebook (Feb 2004 to present - Palo Alto, California) and Harvard University (Computer Science - Psychology, with courses CS182 and CS121). It also lists his high school, Ardsley High School, and Phillips Exeter Academy (Class of 2002). The 'Philosophy' section features a quote: "All children are artists. The problem is how to remain an artist once he grows up." The right sidebar contains a 'You and Mark' section with 3 mutual friends, a 'Sponsored' section with ads for Police Auctions, SF Bucket List, Stay close to your team, and Craft Beer Attorney, and a 'Create an Ad' link.

- How does Facebook use your data?
- Where do you think Facebook can use your data?


# What about Amazon?





English ▾

Friends' Activity **0** Sign Up for Yelp Log In



**Real people. Real reviews.®**

**Search for** (e.g. taco, cheap dinner, Max's)

**Near** (Address, Neighborhood, City, State or Zip)

San Francisco Community Acupuncture

San Francisco Food Bank

San Francisco Symphony

**san francisco restaurants**

San Francisco Movie Tours

San Francisco Marriott Fishermans Wharf

San Francisco Gay Wedding Video

San Francisco Soup Company

San Francisco Test Only Smog

San Francisco CityPASS

Talk Events

Member Search

**Try Yelp in...**

Are You Looking For **Yelp Berkeley**

Amsterdam

Atlanta

Austin

Berlin

Boston

Chicago

Dallas

Denver

Detroit

Dublin

Honolulu

Houston


**Yelp Berkeley**

Yelp is the fun and easy way to find and talk about great (and not so great) local businesses

**Best of yelp**

**Restaurants**

4505 reviewed



1. La Bedaine
2. Kingston 11 Cuisine
3. Vital Vittles
4. Cheese Board Pizza
5. Emilia's Pizzeria


...see more »

**Nightlife**

881 reviewed

**Shopping**

4852 reviewed



1. Chestnut & Vine Floral...
2. Waterside Workshops
3. UniFormal & UniEleganza Tuxedo...
4. Lee's Florist & Nursery
5. Supple Integrative Skin Care

...see more »

**Beauty and Spas**

2566 reviewed

**Browse by Category**


- 🍴 Restaurants
- 🍴 Food
- 🍷 Nightlife
- 🛒 Shopping
- 💆 Beauty and Spas
- 🎨 Arts & Entertainment
- 📅 Event Planning & Services
- 🏥 Health and Medical
- 🏠 Active Life
- ✈ Hotels & Travel
- 🚗 Automotive
- 🏠 Home Services
- 🚗 Local Services
- 🌟 Local Flavor
- 🐾 Pets
- 🎓 Education

**Get the Yelp app on your mobile phone**

It's free and helps you find great, local businesses on the go!

**Review of the Day**

Voted by our members!

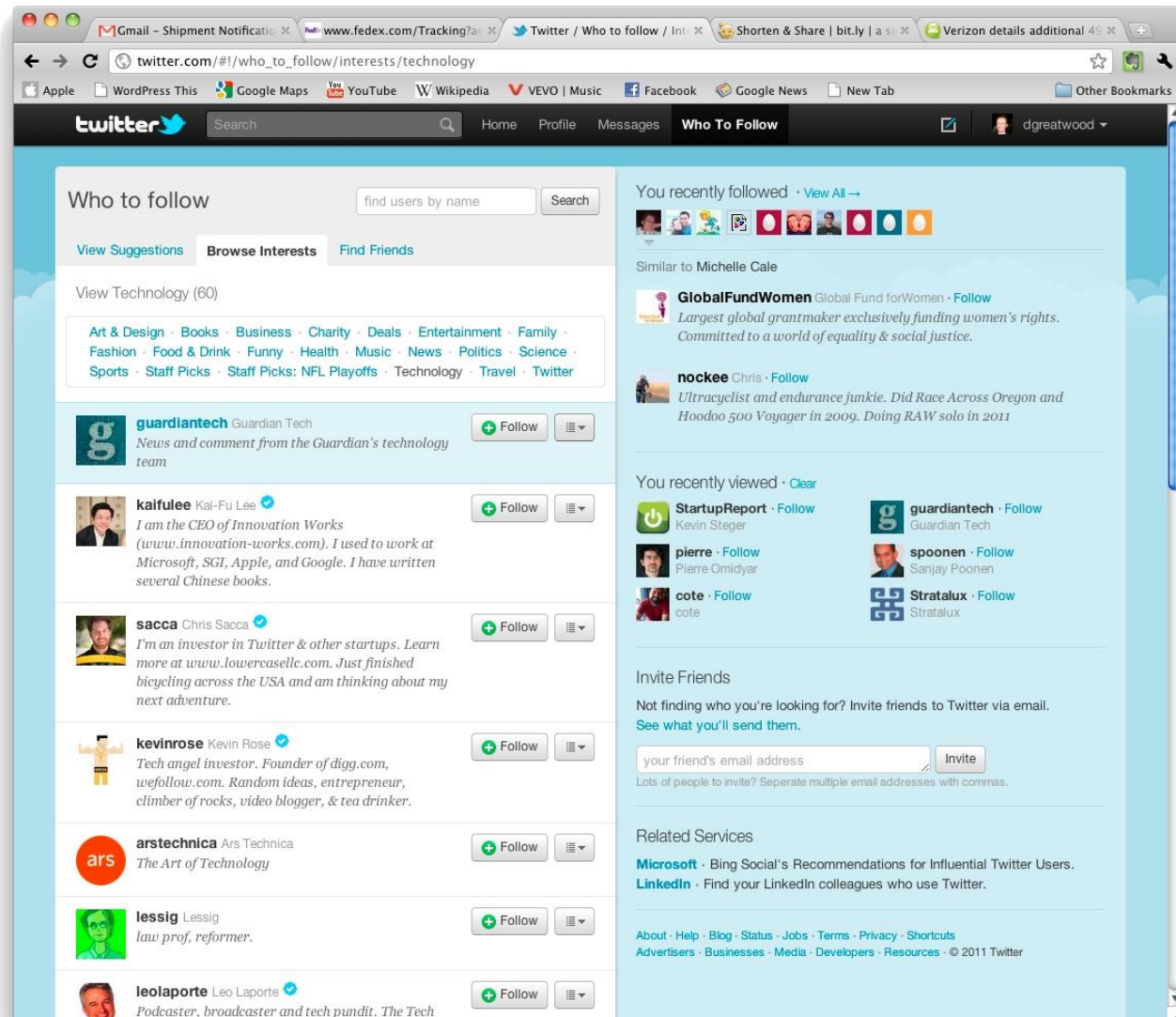


★★★★★

Dopo

\*Cozy, romantic little neighborhood spot: Nice location right across from good ol' Fenton's. No reservations at this little place. so either come a little

# Or Twitter?



## Social Media Mining

- Line number:
  - CIS 700
- Priority for students that have taken CIS 787
  - See blackboard for more information / Follow the procedure to obtain my signature
- Classroom and Hours:
  - CST 3-216, MW 3:45 – 5:05 PM
- Blackboard:
  - Everything (slides/homeworks/projects/etc.)

- **Instructor:**
  - Reza Zafarani (rzafaran@syr.edu)
- **Office Hours (Starting Jan 28th):**
  - Thursday 9:00 - 10:00am, CST 4-279
  - Other times: by appointment only



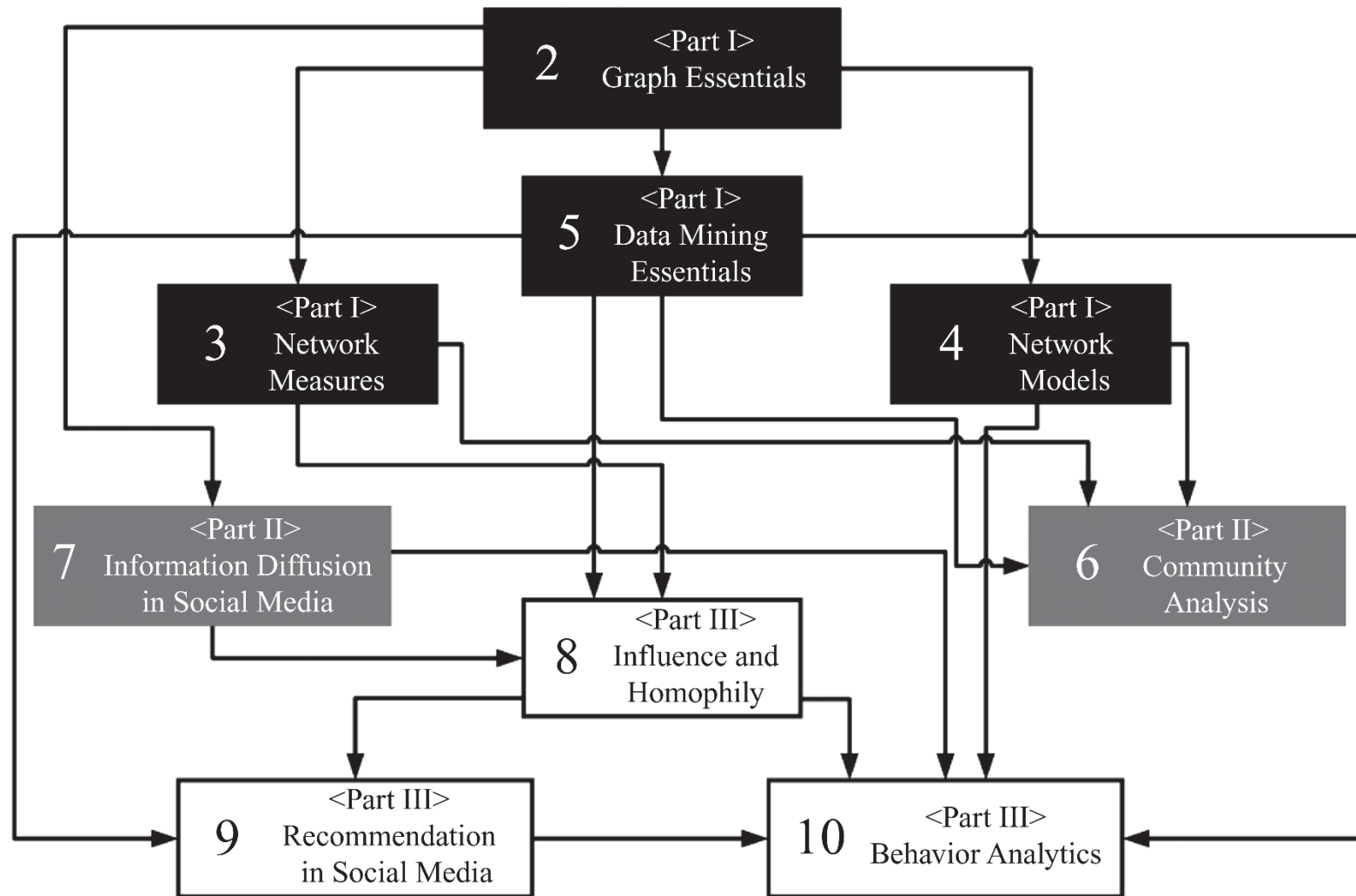
# Objectives of Our Course

- Understand social aspects of the Web
  - Social Theories + Social media + Mining
  - Learn to collect, clean, and represent social media data
  - How to measure important properties of social media and simulate social media models
  - Find and analyze communities in social media
  - Understand how information propagates in social media
  - Understanding friendships in social media, perform recommendations, and analyze behavior
- Study or ask interesting research issues
  - e.g., start-up ideas / research challenges
- Learn representative algorithms and tools

# Course Information

- Prerequisites:
  - CIS 787 – Analytical Data Mining
  - Data Structures and Algorithms
    - Search/Sort algorithms
    - Graphs
    - Graph Algorithms (Traversal, MST, shortest-paths)
    - Time/Space Complexity
  - Programming Skills: Java, basic understanding of MATLAB is a plus
    - E.g., Being able to crawl a website with Java
    - E.g., Computing eigenvalues of a matrix with MATLAB
  - Basic knowledge of probability, statistics, calculus, and linear algebra
    - Expectation, variance, standard deviation,
    - Eigenvalue computation, determinants, characteristic equation
    - Basic differentiation, integration, and differential equations

# Overview – Dependency Graph



The weekly schedule  
will be available on  
blackboard

# Course Workload and Evaluation

- **A lot** of work is expected from you! Think twice if it does not fit your schedule or match your expectation.
  - Lectures
    - Experienced researchers or practitioners may be invited as guest instructors for specific topics.
  - Homework assignments (15%) – 4 HWs
    - Conceptual; deep thinking required
  - Projects (20%)
    - Two projects
      - first project is individual, second is group project
  - 3 Exams (50%) – 3/6: highest, 2/6: medium, 1/6: lowest
  - Quizzes (15%) – 9-10 quizzes; after chapters/topics
  - Late penalty:
    - Exponential penalty: -50% first day, -75% second day, no points on or after the third day.



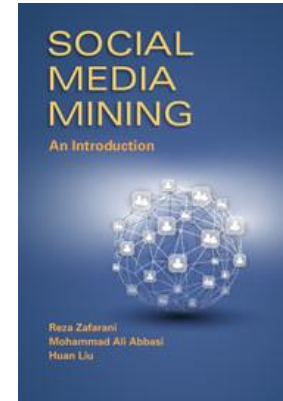
# Academic Integrity

- Please see:
  - [http://supolicies.syr.edu/ethics/acad\\_integrity.htm](http://supolicies.syr.edu/ethics/acad_integrity.htm)
- You are encouraged to form groups to solve problems and coding assignments; however, when writing, write in your own words and provide your own solutions.

## Primary Reference:

Social Media Mining, Reza Zafarani,  
Mohammad-Ali Abbasi, Huan Liu,  
Cambridge University Press 2014,

- Available at  
<http://socialmediamining.info> or  
amazon.com



# Communication Channels and Schedule

## Me → You

- Announcements are made regularly on Blackboard
  - Check blackboard regularly
- Emails will be sent out on a need basis

## You → me

- Office hours / Email (I don't check my voicemail)

## Many ↔ Many

- Q & A: You can ask questions from the instructor and/or other students on blackboard.
  - After-class Blog (i.e., Discussion Board)

# Feedback

- A class survey today
- One in the middle of the semester
- At the end of each of classes:

Topic	Yes	No
Topic 1	X	
Topic 2		X
Topic 3	X	
Topic 4		X
"the equation on slide 5"		X

- All surveys are anonymous

# Social Media



# Definition

Social Media is the use of electronic and Internet tools for the purpose of sharing and discussing information and experiences with other human beings in more efficient ways.

# Social Media Landscape 2015



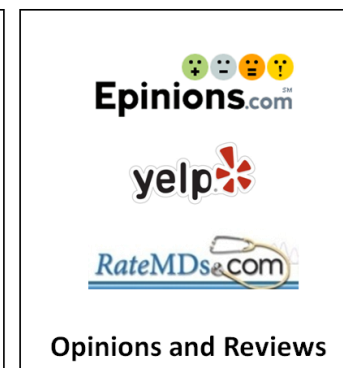
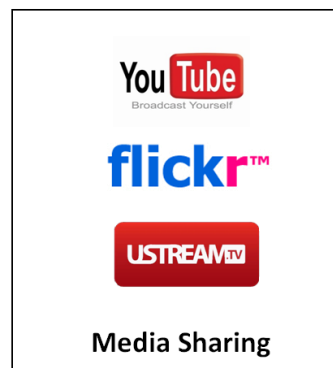
FredCavazza.net

# Social Media: Examples

- A wiki article
- Web reviews and ratings of a popular pizza place in your city
  - E.g., Yelp.com
- An online social network of your professional contacts
  - E.g., Facebook.com, LinkedIn.com
- An iPhone application that informs you where parking is likely available
  - FasPark

# Types of Social Media

- Online Social Networking
- Publishing
  - Blogging
  - Wiki
- Micro blogging
- Social News
- Social Bookmarking
- Media Sharing
  - Video Sharing
  - Photo Sharing
  - Podcast Sharing
- Opinion, Review, and Ratings Websites
- Answers
- Entertainment



# Online Social Networking

Online Social Networks are web-based services that allow individuals and communities to connect with real world friends and acquaintances online

- Interactions
  - Friendship interaction
    - Friends, like, comments, ...
  - Media Sharing
  - Sending and receiving messages

- Examples
  - Facebook.com
  - MySpace.com
  - Bebo.com
  - Orkut.com

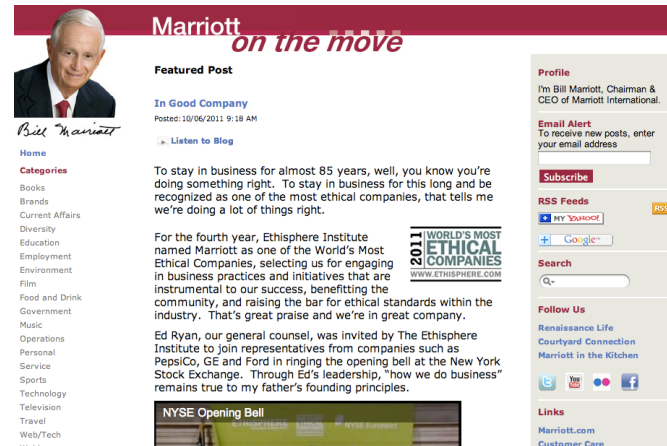




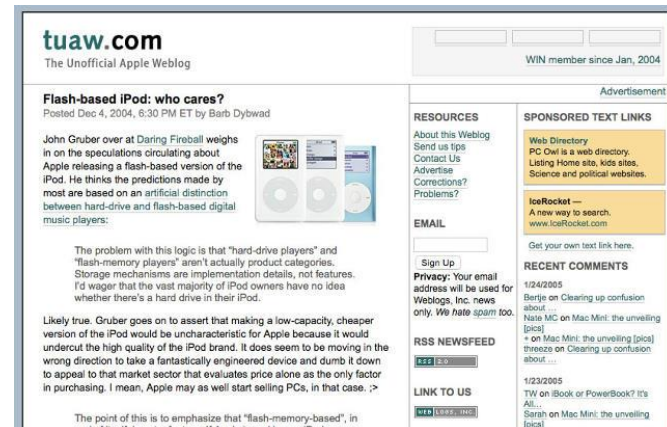
# Blogging

A blog is a journal-like website for users, a.k.a. bloggers, to contribute textual and multimedia content, arranged in reverse chronological order

- Maintained both individually or by a community
  - See a tutorial at KDD  
[http://videolectures.net/kdd08\\_liu\\_briat/](http://videolectures.net/kdd08_liu_briat/)
- Usages:
  - Sharing information and opinions with friends and strangers
  - Disseminating subject-specific content
  - Who is the influential  
[http://videolectures.net/wsdm08\\_agarwal\\_iib/](http://videolectures.net/wsdm08_agarwal_iib/)



The screenshot shows a blog post by Bill Marriott, Chairman & CEO of Marriott International. The post is titled "In Good Company" and is dated 10/06/2011 9:18 AM. The post content discusses Marriott's recognition as one of the World's Most Ethical Companies for the fourth year, selected by Ethisphere Institute. It mentions that Marriott is doing something right and staying in business for almost 85 years. The post also mentions that Marriott is being recognized as one of the most ethical companies, that tells me we're doing a lot of things right. The post is categorized under "Featured Post" and "In Good Company". There are links to "Listen to Blog" and "Subscribe". The sidebar includes a "Profile" section, an "Email Alert" section, an "RSS Feeds" section, a "Search" section, a "Follow Us" section, and a "Links" section.



The screenshot shows a blog post titled "Flash-based iPod: who cares?" on the website tuaw.com. The post is dated Dec 4, 2004, 6:30 PM ET by Barb Dybwad. The post content discusses the problem with the iPod being a "hard-drive player" and "flash-memory players" not being actual product categories. It mentions that Apple releasing a flash-based version of the iPod would be uncharacteristic for Apple because it would undercut the high quality of the iPod brand. The post also mentions that Apple may be well start selling PCs, in that case. The post is categorized under "Flash-based iPod: who cares?". There are links to "About this Weblog", "Send us tips", "Contact Us", "Advertise", "Corrections?", "Problems?", "EMAIL", "Sign Up", "Privacy", "RSS NEWSFEED", "LINK TO US", "SPONSORED TEXT LINKS", "RECENT COMMENTS", and "WIN member since Jan, 2004".

# Microblogging

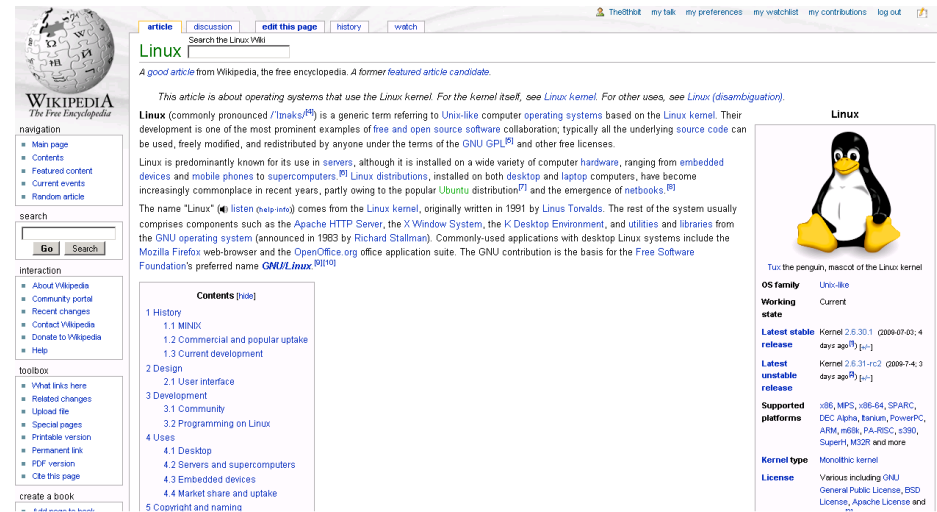
Microblogging can be considered as a counterpart to blogging, but with limited content

- Usage
  - communication medium
  - social interaction
  - citizen journalism
- Service Providers:
  - Twitter
  - Google buzz



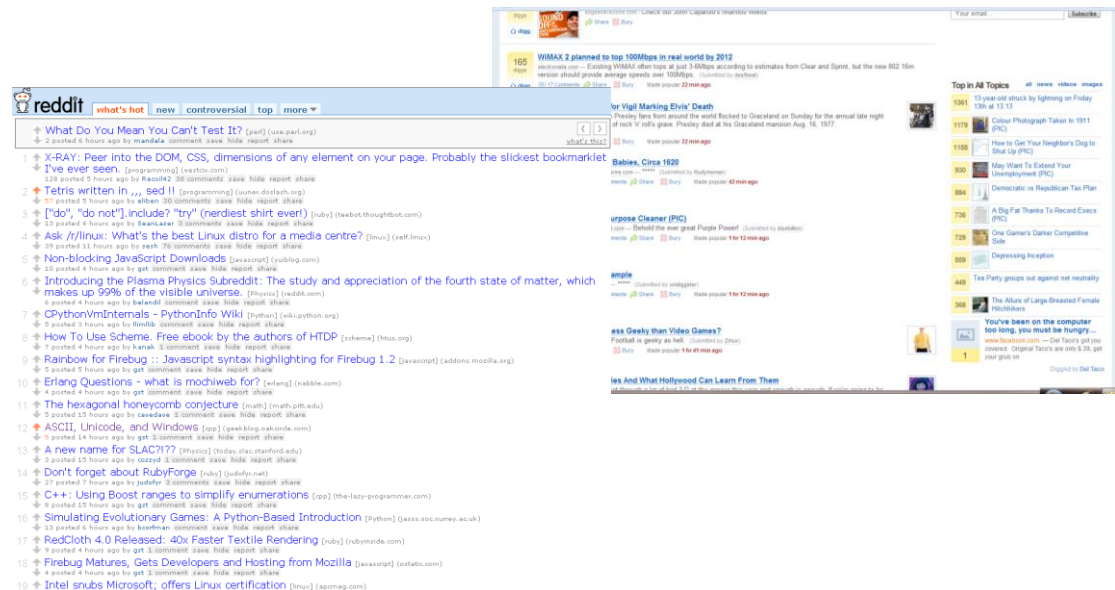
A wiki is a collaborative editing environment that allows users to develop Web pages using a simplified markup language

- Wikipedia allows interested individuals to collaboratively develop articles on a variety of subjects.
- Using the wisdom of crowds effectively, it has become a comprehensive repository of information useful to a variety of individuals



Social News refers to the sharing and selection of news stories and articles by a community of users.

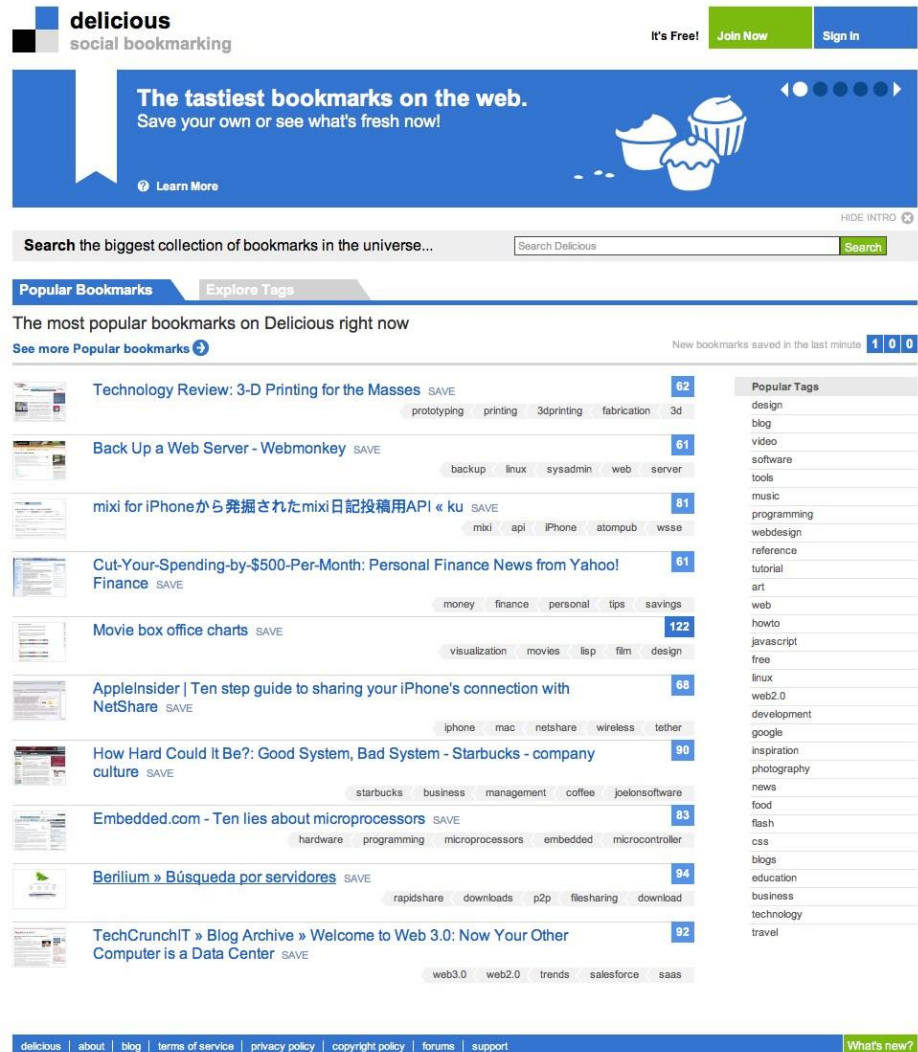
- Users can share articles that they believe would interest the community
- Samples:
  - Digg.com
  - Slashdot
  - Fark
  - Reddit



# Social Bookmarking

Social Bookmarking sites allow users to bookmark web content for storage, organization and sharing.

- These bookmarks can be tagged with metadata to categorize and provide context to the shared content, allowing users to organize information making it easy to search and identify relevant information.
- Samples
  - Delicious.com
  - StumbleUpon.com



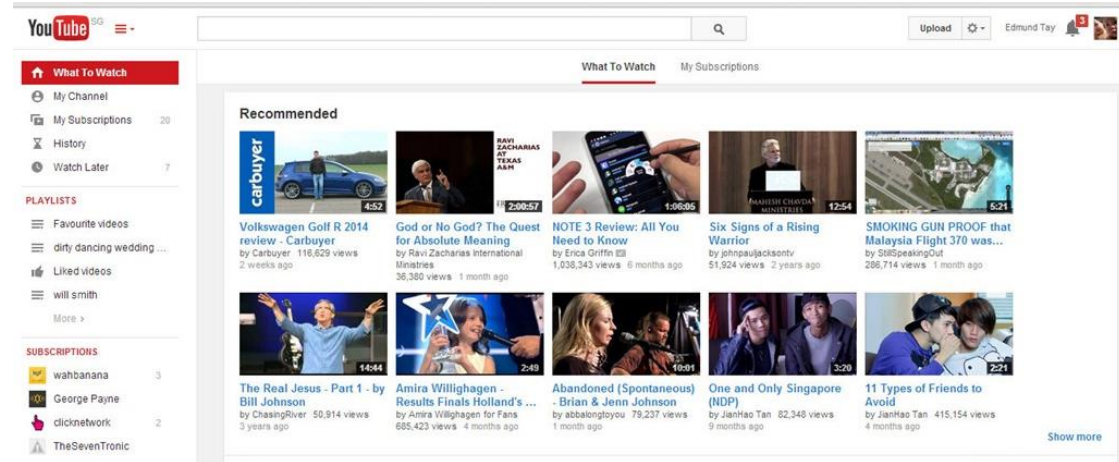


# Media Sharing

Media sharing is an umbrella term that refers to the sharing of a variety of media on the web.

Users share such multimedia content of possible interest to others

- Samples:
  - Video Sharing:
    - YouTube.com
  - Photo Sharing:
    - Flickr.com, picasa.com
  - Document Sharing:
    - Scribd.com, Slideshare.com
  - Livecasting:
    - Justin.tv, Ustream.com



# Opinion, Review, and Ratings Websites

Opinion, review, and ratings websites are websites whose primary function is to collect and publish user-submitted content in the form of subjective commentary on existing products, services, entertainment, businesses, places, etc. Some commercial sites may serve a secondary purpose as review sites by publishing product reviews submitted by customers.

- Examples
  - Cnet.com
  - Epinions.com
  - yelp.com
  - tripadvisor.com

The screenshot displays the Yelp website interface. At the top, there is a search bar with the text "Search for (e.g. taco, cheap dinner, Max's)" and a location dropdown set to "sf, ca". Below the search bar is a navigation menu with links: "Welcome", "About Me", "Write a Review", "Find Friends", "Messaging", "Talk", "Events", and "Member Search". The main content area shows a breadcrumb trail: "Tartine Bakery > Menu > Breakfast Pastries > Croissant". The "Croissant" item is highlighted with a price of "\$3.85" and three photos. Below the photos, there is a section for reviews. A review by Stephanie S. from Loma Linda, CA, dated 10/23/2012, is shown. The review text reads: "This was our first stop from the airport and we were starving! The line was long, but it went pretty fast. This was our first time here and we couldn't decide what to order. We tried the morning bun, chocolate and almond croissant, bread pudding, & the chocolate eclair. Everything was delicious, but the morning bun was soo amazing. I loved the hints of citrus and the flakiness of the bun. I made my hubby go back & buy me another one to save for later. Oh, Tartine! I wish you were also located in So. Cal." Below the review, there are buttons for "Write a Review", "Add a photo", "Compliment", "Send Message", and "Follow This Review". On the right side, there is a "Menu for Tartine Bakery" section with a list of items and their prices: "Croissant" (\$3.85), "Frangipane Croissant" (\$4.50), "Double Pain Au Chocolat" (\$4.50), "Morning Buns" (\$3.85), "Buttermilk Scones" (\$3.25), "Tea Cake" (\$3.75), "Bread Pudding" (Price details), "Pain Au Jambon" (\$4.95), "Gougere" (\$3.50), "Cake Aux Olives" (\$4.95), "Quiche" (Price details), and "Muesli" (Price details). Each item has a small icon, a number of reviews, and a number of photos.

# Socially-Provided Answers

In these sites, users who require certain guidance, advice or knowledge can ask questions. Other users from the community can answer these questions based on knowledge acquired from previous experiences, personal opinions or from relevant research.

- Unlike review and opinion sites, which contain self-motivated contribution of opinions, answer sites contain knowledge shared in response to a specific query.
- Samples:
  - WikiAnswers, Yahoo Answers, Quora

The screenshot displays a Quora interface with a search bar at the top containing the text "Search Google Analytics Questions and Topics" and an "Add Question" button. Below the search bar, there are two question entries. The first question is "What percentage of visits would Omniture / Google Analytics / Coremetrics etc miss?" with the answer "Assuming client-side integration, compared with the numbers from the web servers and proxy logs." and 0 answers. The second question is "How can I track Pinterest in Google Analytics?" with the answer "Their Javascript pinit.js file (http://assets.pinterest.com/js/p... c) doesn't seem to add any callbacks, so the best you can do is track clicks on the 'Pin It' button in Goo..." and 2 answers. To the right of the questions, there is a section titled "Share Topic · Invite People" with links to Twitter, Facebook, and Quora. Below that, there is a "Top Answerers" section listing Mike Sullivan, Ozberk Olcer, Shay Sharon, AJ Kohn, and Christopher O'Donnell. At the bottom right, there is a "Followed by 5455 People" section showing a grid of user avatars.

# Main Characteristics

- **Participation**
  - social media encourages contributions and feedback from everyone who is interested. It blurs the line between media and audience.
- **Openness**
  - most social media services are open to feedback and participation. They encourage voting, comments and the sharing of information. There are rarely any barriers to accessing and making use of content – password-protected content is frowned on.
- **Conversation**
  - whereas traditional media is about “broadcast” (content transmitted or distributed to an audience) social media is better seen as a two-way conversation.
- **Community**
  - social media allows communities to form quickly and communicate effectively. Communities share common interests, such as a love of photography, a political issue or a favorite TV show.
- **Connectedness**
  - Most kinds of social media thrive on their connectedness, making use of links to other sites, resources and people.

***Social Media Mining*** is the process of representing, analyzing, and extracting meaningful patterns from social media data

# Social Media Mining Challenges

## 1. Big Data Paradox

1. Social media data is big, yet not evenly distributed.
2. Often little data is available for an individual

## 2. Obtaining Sufficient Samples

1. Are our samples reliable representatives of the full data?

## 3. Noise Removal Fallacy

1. Too much removal makes data more sparse
2. Noise definition is relative and complicated and is task-dependent

## 4. Evaluation Dilemma

1. When there is no ground truth, how can you evaluate?

# TODO Items

- 3 To-do items for you:
  - Follow the process on blackboard to obtain my permission
  - Familiarize yourself with blackboard
  - Take the course survey