

# Web Science: Collective Intelligence & Recommender Systems (Part 1 - Collective Intelligence)

CS 432/532

Old Dominion University

*Permission has been granted to use these slides from Frank McCown, Michael L. Nelson, Alexander Nwala, Michele C. Weigle*



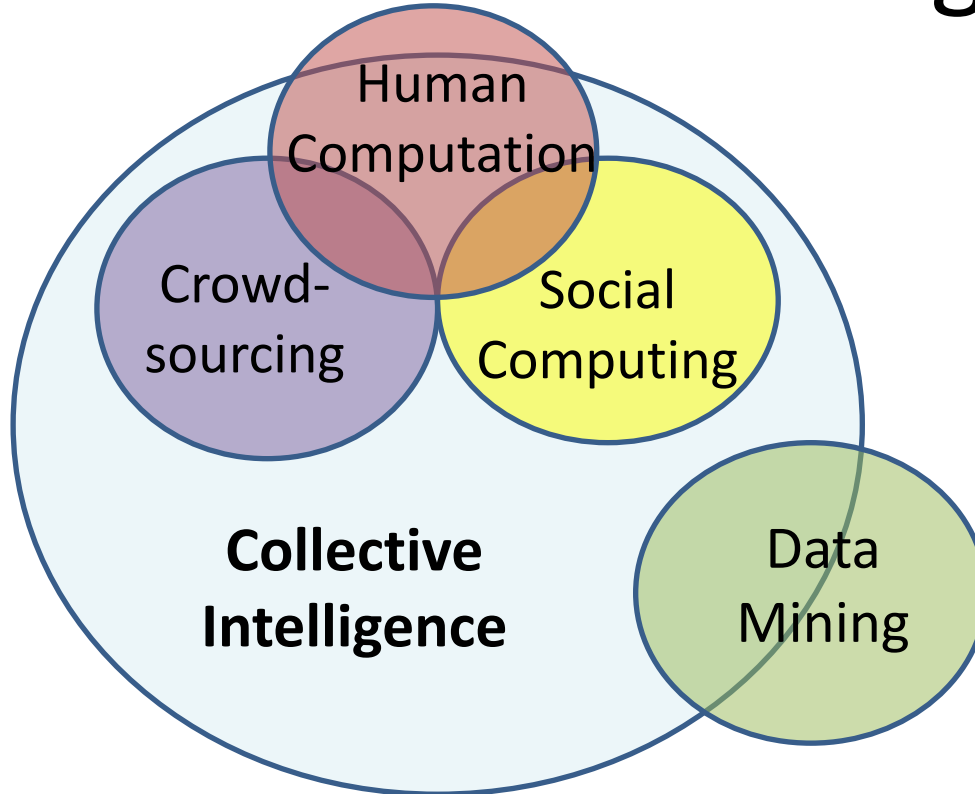
This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License](https://creativecommons.org/licenses/by-nc-sa/3.0/)

## Main reference:

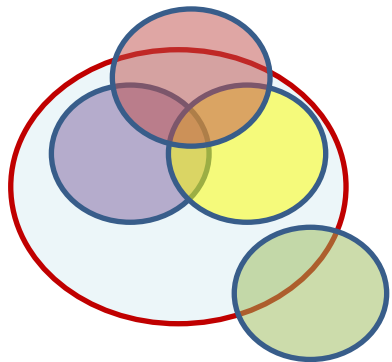
Ch 1 from [Programming Collective Intelligence](#) by Toby Segaran

*(abbreviated as PCI)*

# What is Collective Intelligence?



[Quinn and Bederson](#), "Human computation: A survey and taxonomy of a growing field", CHI 2011



# Collective Intelligence

"groups of individuals doing things collectively that seem intelligent."

[Thomas Malone](#), Director of MIT Center for Collective Intelligence, 2006

"when technologists use this phrase they usually mean the combining of behavior, preferences, or ideas of a group of people to create novel insights."

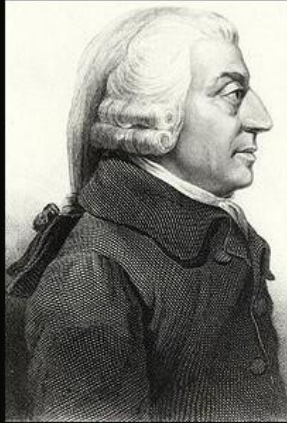
Segaran, *Programming Collective Intelligence (PCI)*, p. 2



img source: <https://www.shutterstock.com/image-photo/guess-how-many-jelly-beans-mason-121109632>

Crowd-sourced strategies: [How to Count Jelly Beans in a Jar](#), [How to win a guess the number of jelly beans in a jar contest](#)

# Collective Intelligence in Classical Economic Theory



It is not from the benevolence of the butcher, the brewer, or the baker that we expect our dinner, but from their regard to their own interest.

(Adam Smith)

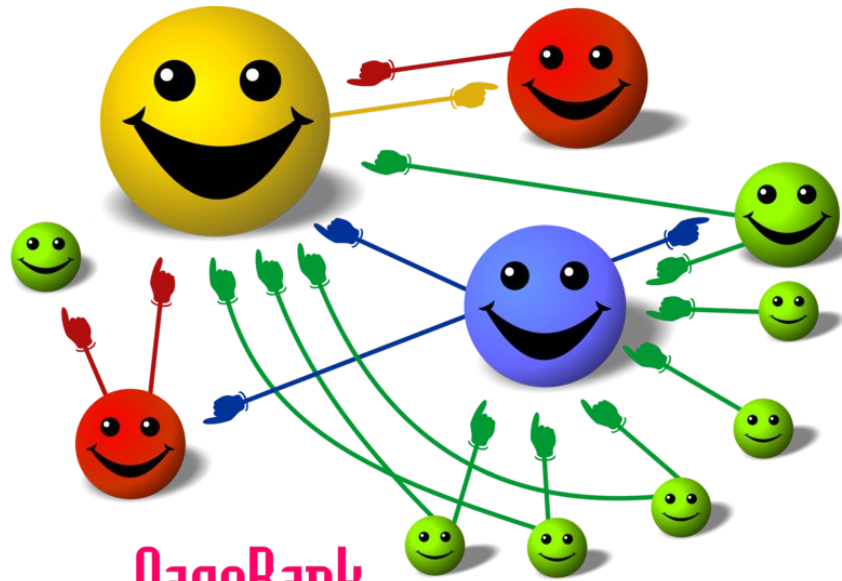
izquotes.com

[Invisible hand](#) (Wikipedia)  
[The Wealth of Nations](#) (Wikipedia)

# Google

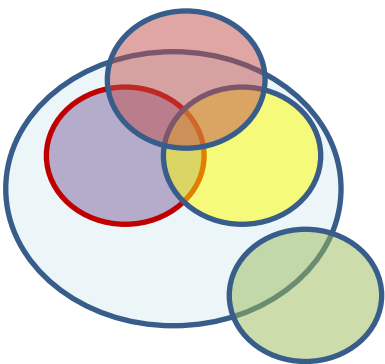
Google Search

I'm Feeling Lucky



## PageRank

Img source: <http://organicseoexpert.org/wp-content/uploads/2012/03/pagerank.png>



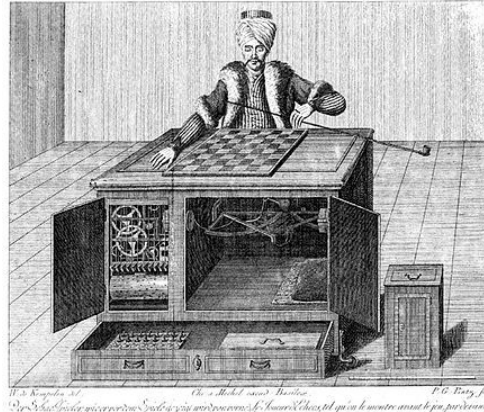
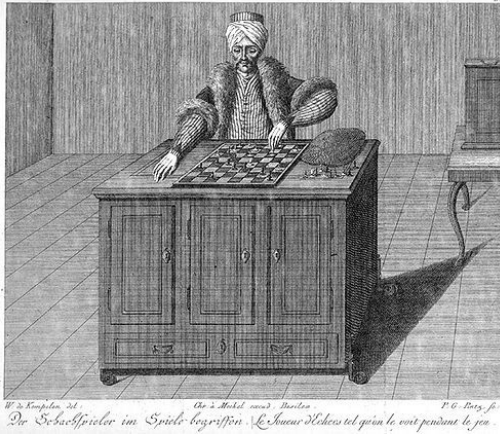
# Crowdsourcing

"the act of taking a job traditionally performed by a designated agent (usually an employee) and outsourcing it to an undefined, generally large group of people in the form of an open call."

Jeff Howe, ["The Rise of Crowdsourcing"](#), *Wired*, Jun 2006



# Mechanical Turk



## The Turk (Wikipedia)

All HITS | HITS Available To You | HITS Assigned To You

Find  containing  that pay at least \$  ☐ for which you are qualified ☐ require Master Qualification

## All HITS

1-10 of 1285 Results

Sort by:

[Show all details](#) | [Hide all details](#)

1 2 3 4 5 > [Next](#) >> [Last](#)

3 questions about your city UNDER 230,000 population only = \$0.17 bonus!\*\*\* - qualification instantly granted (no wait) [View a HIT in this group](#)

**Requester:** [WSOVC.COM](#)

**HIT Expiration Date:** Jul 11, 2012 (3 weeks 6 days)

**Reward:** \$0.00

**Time Allotted:** 2 hours

**HITS Available:** 23327

Help Us Find a URL's Page Ranking on Google (CA) [View a HIT in this group](#)

**Requester:** [CrowdSource](#)

**HIT Expiration Date:** Jun 13, 2013 (52 weeks)

**Reward:** \$0.10

**Time Allotted:** 60 minutes

**HITS Available:** 14999

Give Your Opinion - Simple and Quick! (US) [View a HIT in this group](#)

**Requester:** [CrowdSource](#)

**HIT Expiration Date:** Jun 13, 2013 (52 weeks)

**Reward:** \$0.16

**Time Allotted:** 32 minutes

**HITS Available:** 14978

Help Us Find a URL's Page Ranking on Google (US) [View a HIT in this group](#)

**Requester:** [CrowdSource](#)

**HIT Expiration Date:** Jun 13, 2013 (52 weeks)

**Reward:** \$0.10

**Time Allotted:** 60 minutes

**HITS Available:** 14975

Inv\_B\_2 [View a HIT in this group](#)

**Requester:** [rohzi0d](#)

**HIT Expiration Date:** Jun 28, 2012 (2 weeks 1 day)

**Reward:** \$0.00

**Time Allotted:** 48 minutes

**HITS Available:** 8375

Choose the best category for this link [View a HIT in this group](#)

**Requester:** [Larry Fitzgibbon](#)

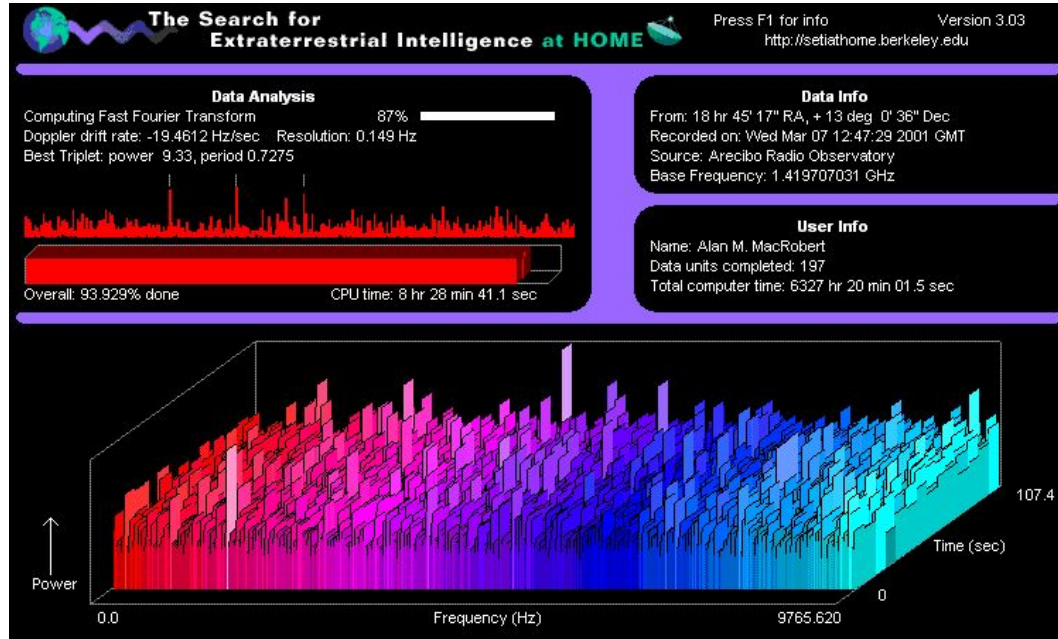
**HIT Expiration Date:** Jun 15, 2012 (2 days 14 hours)

**Reward:** \$0.05

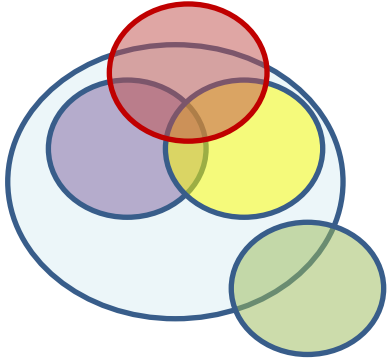
**Time Allotted:** 60 minutes

**HITS Available:** 5000

# SETI@home Harvests Cycles From Idle Computers



[SETI@home](http://setiathome.berkeley.edu)



# Human Computation

"a paradigm for utilizing human processing power to solve problems that computers cannot yet solve."

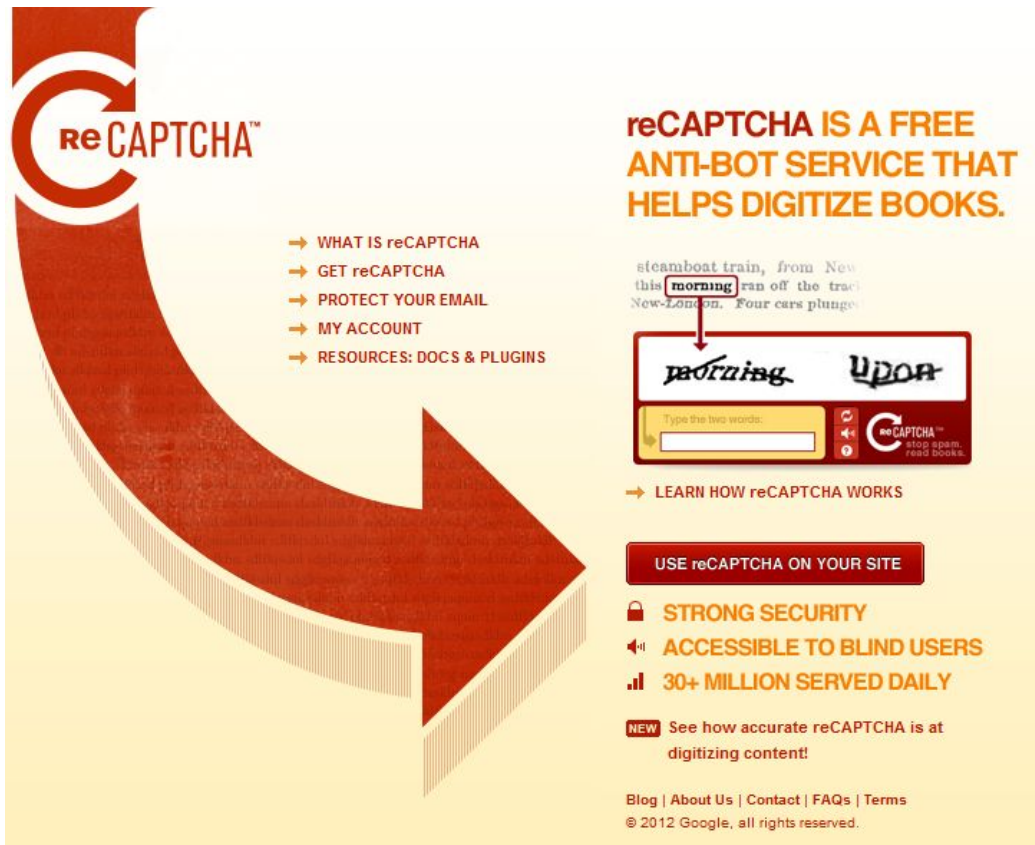
Luis van Ahn, Doctoral Dissertation at Carnegie Mellon, 2005

"Human computation:

The problems fit the general paradigm of computation, and as such might *someday* be solvable by computers.

The human participation is directed by the computational system or process."

Quinn and Bedderson, CHI 2011



The image shows the reCAPTCHA website landing page. On the left, a large red arrow curves downwards. At the top left is the reCAPTCHA logo. Below it, a list of links: WHAT IS reCAPTCHA, GET reCAPTCHA, PROTECT YOUR EMAIL, MY ACCOUNT, and RESOURCES: DOCS & PLUGINS. The main heading reads 'reCAPTCHA IS A FREE ANTI-BOT SERVICE THAT HELPS DIGITIZE BOOKS.' Below this, a snippet of text from a digitized book is shown: 'steamboat train, from New this morning ran off the track New-London. Four cars plunges'. A red box highlights the word 'morning'. Below the text is a reCAPTCHA interface showing the words 'morning' and 'upon' in a stylized font. A text input field is labeled 'Type the two words:'. To the right of the input field is a reCAPTCHA logo with the text 'stop spam. read books.' Below the interface, a link 'LEARN HOW reCAPTCHA WORKS' is shown. A red button says 'USE reCAPTCHA ON YOUR SITE'. Below this, three features are listed: 'STRONG SECURITY', 'ACCESSIBLE TO BLIND USERS', and '30+ MILLION SERVED DAILY'. A 'NEW' badge is next to the text 'See how accurate reCAPTCHA is at digitizing content!'. At the bottom, there are links for 'Blog | About Us | Contact | FAQs | Terms' and a copyright notice '© 2012 Google, all rights reserved.'

reCAPTCHA™

- WHAT IS reCAPTCHA
- GET reCAPTCHA
- PROTECT YOUR EMAIL
- MY ACCOUNT
- RESOURCES: DOCS & PLUGINS

reCAPTCHA IS A FREE ANTI-BOT SERVICE THAT HELPS DIGITIZE BOOKS.

steamboat train, from New this **morning** ran off the track New-London. Four cars plunges

morning upon

Type the two words:

reCAPTCHA™ stop spam. read books.

→ LEARN HOW reCAPTCHA WORKS

USE reCAPTCHA ON YOUR SITE

- 🔒 STRONG SECURITY
- 🔊 ACCESSIBLE TO BLIND USERS
- 📊 30+ MILLION SERVED DAILY

**NEW** See how accurate reCAPTCHA is at digitizing content!

[Blog](#) | [About Us](#) | [Contact](#) | [FAQs](#) | [Terms](#)

© 2012 Google, all rights reserved.

[reCAPTCHA v3](#), [Google's New Street View Image Recognition Algorithm Can Beat Most CAPTCHAs](#)  
Goodfellow et al., ["Multi-digit Number Recognition from Street View Imagery using Deep Convolutional Neural Networks"](#), 2014





[How Google Cracked House Number Identification in Street View](#)



## Tasks (XKCD)

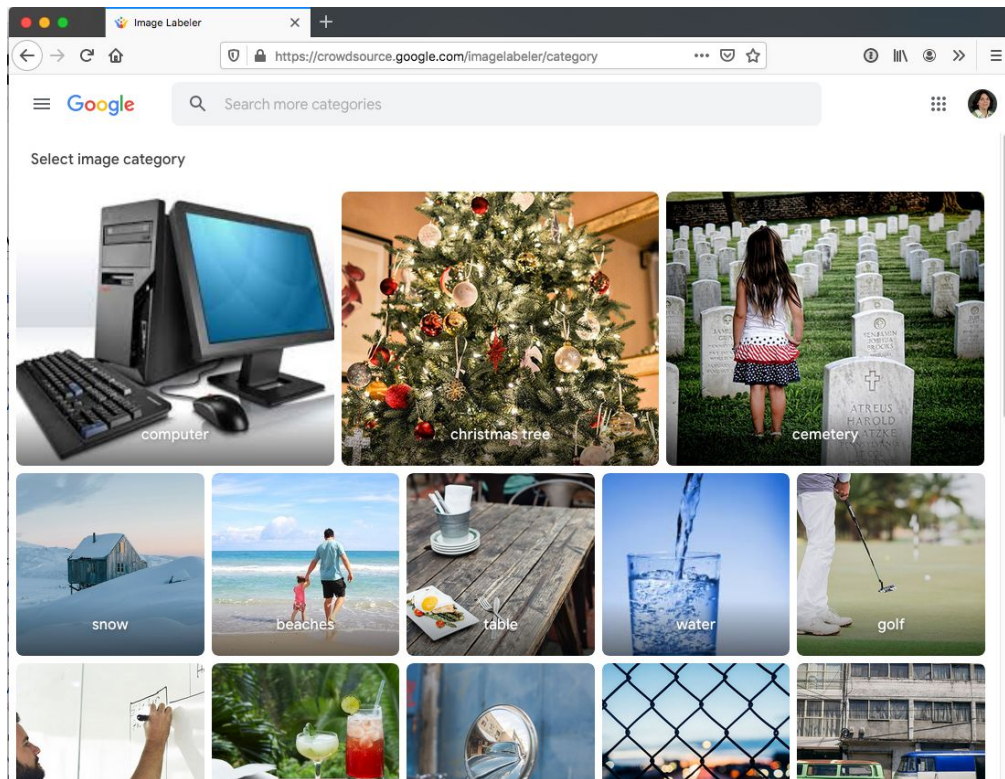
In the 60s, Marvin Minsky assigned a couple of undergrads to spend the summer programming a computer to use a camera to identify objects in a scene. He figured they'd have the problem solved by the end of the summer. Half a century later, we're still working on it.



[ESP game](#) (Wikipedia)



# Google licensed ESP for ImageLabeler



[ImageLabeler](#) (Google)

# FoldIt

Shake sidechains to improve the protein.  
Hotkey: S

Shake Sidechains Wiggle Backbone Clear Locks and Bands Reset Puzzle Mouse Help

▲ Actions ► History ► View ► File

Pull Tool

Rank: 17 Score: 9092  
48: Pro Peptide

▼ Group Competition

#	Group Name	Score
1	The Lone Folder	9388
2	Street Smarts	9367
3	Illinois	9303
4	Berkeley	9255

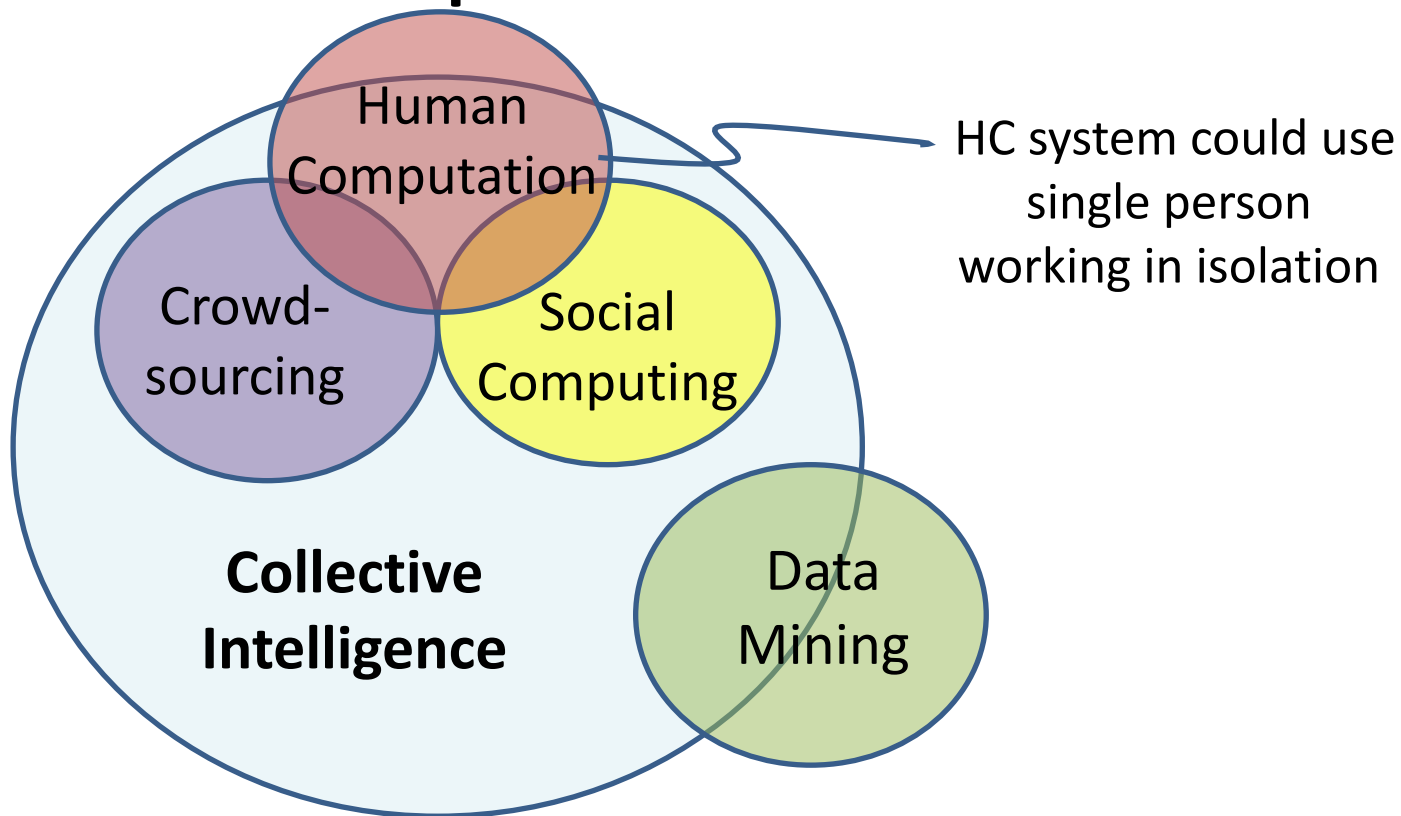
▼ Player Competition

16	psen	-	9098
17	kathleen	9092	9092
18	versat82	-	9091
19	darktorres	-	9081
20	carrico	9032	9066
21	mbjorkegren	-	9048
22	sslickerson	-	9038

► Chat

[Foldit: Solve Puzzles for Science](#)

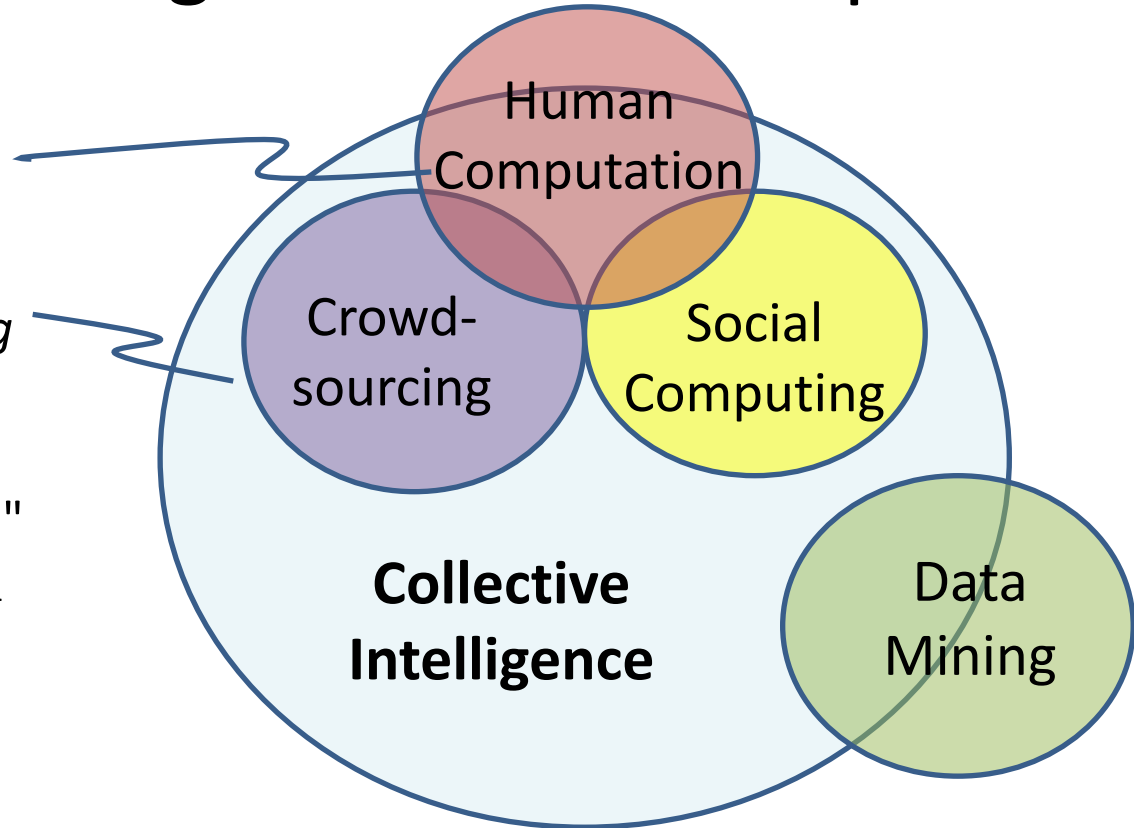
# CI not superset of HC



# Crowdsourcing vs. Human Computation

"Whereas *human computation* replaces computers with humans, *crowdsourcing* replaces traditional human workers with members of the public."

Quinn and Bedderson, CHI 2011



# DARPA Network Challenge



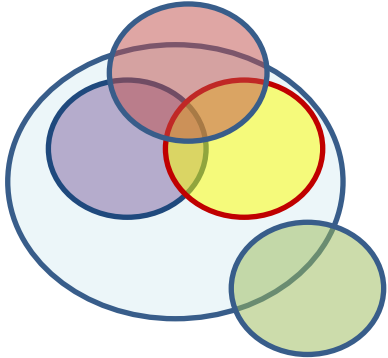
[DARPA Network Challenge](#) (Wikipedia)

**Challenge:** Find 10 red weather balloons set up by DARPA, from 10am - 5pm EST, Dec 5, 2009. They were prepared to run the challenge for 1 week.

MIT placed first (of 10 teams). They found all 10 in less 9 hours.

**MIT's Strategy:** invitation-based network (with 4 initial members), finder gets \$2000, with everyone in the invitation network getting 0.5 of the previous amount, with the remainder donated to charity.

MIT site: [MIT Red Balloon Challenge \(archived\)](#)



# Social Computing

"applications and services that facilitate collective action and social interaction online with rich exchange of multimedia information and evolution of aggregate knowledge."

Parameswaran and Whinston, *Social Computing: An Overview*, CAIS 19:37, 2007





WIKIPEDIA  
The Free Encyclopedia

[Main page](#)  
[Contents](#)  
[Featured content](#)  
[Current events](#)  
[Random article](#)  
[Donate to Wikipedia](#)

Interaction  
[Help](#)  
[About Wikipedia](#)  
[Community portal](#)  
[Recent changes](#)  
[Contact Wikipedia](#)

[Toolbox](#)  
[Print/export](#)

Languages  
[العربية](#)  
[Česky](#)  
[Deutsch](#)  
[Español](#)  
[Esperanto](#)  
[Français](#)

[Log in](#) [Create account](#)

Article [Talk](#)

[Read](#) [Edit](#) [View history](#)

## Collective intelligence

From Wikipedia, the free encyclopedia

*Not to be confused with [group intelligence](#), [collaborative intelligence](#), or [knowledge sharing](#).*

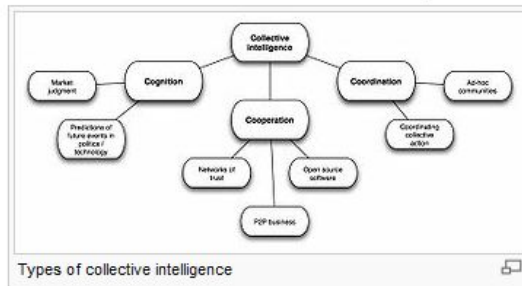


This article **needs attention from an expert on the subject**. Please add a *reason* or a *talk* parameter to this template to explain the issue with the article. [WikiProject Sociology](#) or the [Sociology Portal](#) may be able to help recruit an expert. *(April 2010)*

**Collective intelligence** is a [shared](#) or [group intelligence](#) that emerges from the collaboration and competition of many individuals and appears in [consensus decision making](#) in bacteria<sup>[*[clarification needed](#)*]</sup>, animals, humans and computer networks. The term appears in [sociobiology](#), [political science](#) and in context of [mass peer review](#) and [crowdsourcing](#) applications. This broader definition<sup>[*[clarification needed](#)*]</sup> involves consensus, social capital and formalisms such as [voting systems](#), [social media](#) and other means of quantifying mass activity. *Everything from a [political party](#) to a [public wiki](#) can reasonably be described as this loose form of collective intelligence.*

More narrowly, it can be understood as an emergent property between people and ways of processing information.<sup>[1]</sup> This notion of collective intelligence is referred to as **Symbiotic intelligence** by Norman Lee Johnson.<sup>[2]</sup> The concept is used in [sociology](#), [business](#), [computer science](#) and mass communications: it also appears in science fiction.

Writers who have influenced the idea of collective intelligence include [Douglas Hofstadter](#) (1979), [Peter Russell](#) (1983), [Tom Atlee](#) (1993), [Pierre Lévy](#) (1994), [Howard Bloom](#) (1995), [Francis Heylighen](#) (1995), [Douglas Engelbart](#), [Cliff Joslyn](#), [Ron Dembo](#), [Gottfried Mayer-Kress](#) (2003).



Types of collective intelligence

[Contents](#) [\[hide\]](#)

[Collective intelligence](#) (Wikipedia)



Questions

Tags

Users

Badges

Unanswered

Ask Question

## All Questions

newest

301 featured

faq

votes

active

unanswered

301

questions

6

votes

1

answer

82 views

**+500** ShareKit with MonoTouch how?

How do you use ShareKit with MonoTouch? The MonoTouch Bindings project on GitHub seems to have bindings for ShareKit but I can't get them to work. I currently have an iPhone application developed ...

[monotouch](#)[sharekit](#)

asked Jun 4 at 10:24

[ShareKit with MonoTouch](#)

31 • 1

1

vote

2

answers

75 views

**+50** Web Application and REST services SSO in tomcat and spring-security

I am using two different web application deployed in the same tomcat instance. One of web application and another one is REST services. When user logged into the web application and calls the REST ...

[tomcat](#)[spring-security](#)[sso](#)

asked Jun 5 at 4:36

[Krishna](#)

666 • 6 • 28

2

votes

3

answers

80 views

**+50** Tomcat fails on heavy load

I'm running a heavy db-based GWT application on a debian VPS using tomcat server 7 (& JRE 1.6). My app contains a lot of java servlets which communicate with MySQL5 database via a tomcat ...

[tomcat](#)[server](#)[crash](#)[load](#)[timeout](#)

asked May 22 at 20:35

[Gerhard](#)

2,008 • 6 • 18

16

votes

4

answers

5k views

**+100** Tagging friends in status updates from Facebook API

I recently came across this blog post that said that it's possible to tag someone in a status update from a Facebook app (= from the API): ...

[php](#)[facebook](#)[api](#)[friends](#)

asked Dec 1 '10 at 17:51

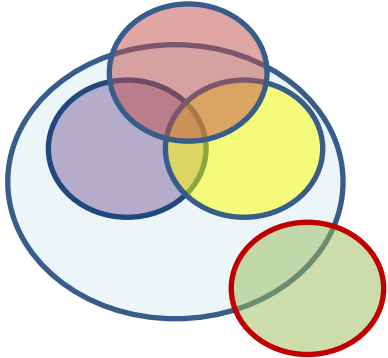
[Vinch](#)

153 • 2 • 8

## Community Bulletin

event **2012 Community Moderator Election** – ends in 6 days**CAREERS 2.0**  
by stackoverflow**CAREERS 2.0**[C# - Web - Remote Work - Software Developer](#)[Dynamic Benchmarking](#)  
Amherst, NH; San Jose, CA[Senior Agile Ruby on Rails Developer](#)

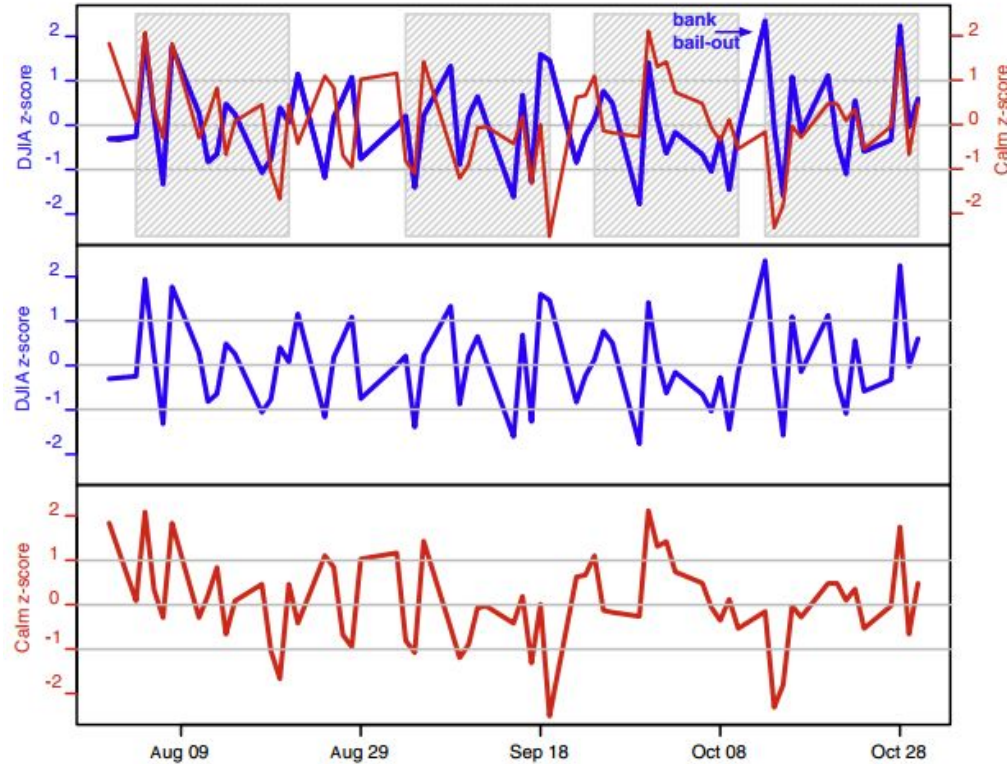




# Data Mining

"the application of specific algorithms  
for extracting patterns from data."

Fayyad, Piatetsky-Shapiro, and Smyth,  
Knowledge Discovery and Data Mining: Towards a  
Unifying Framework, *Proc. KDD*, 1996



Bollen et al., [Twitter mood predicts the stock market](#), 2011

# Web Science: Collective Intelligence & Recommender Systems

(Part 2 - Intro to Recommender Systems)

CS 432/532

Old Dominion University

*Permission has been granted to use these slides from Frank McCown, Michael L. Nelson, Alexander Nwala, Michele C. Weigle*



This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License](https://creativecommons.org/licenses/by-nc-sa/3.0/)

## Main reference this week:

Ch 2 from [Programming Collective Intelligence](#) by Toby Segaran

*(abbreviated as PCI)*

# Recommender Systems

- Recommender (recommendation) systems recommend things to people based on their preferences or past behavior
- Two general approaches:
  - Collaborative filtering
    - You'll like X because other people like you also liked X
  - Content-based (*not covered in these slides*)
    - You'll like X because it is very similar to other things you like

Watch Instantly

Browse DVDs

Your Queue

Movies You'll ♥

## Congratulations! Movies we think You will ♥

Add movies to your Queue, or **Rate** ones you've seen for even better suggestions.

Spider-Man 3



Add



Not Interested

300



Add



Not Interested

The Rundown



Add



Not Interested

Bad Boys II



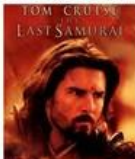
Add



Not Interested

Las Vegas: Season 2  
(6-Disc Series)

The Last Samurai



Star Wars: Episode III

Robot Chicken: Season 3  
(2-Disc Series)

## Just For Today

[Browse Recommended](#)

## Recommendations

[Amazon Instant Video](#)  
[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](#)  
[Industrial & Scientific](#)  
[Jewelry](#)  
[Kitchen & Dining](#)  
[MP3 Downloads](#)  
[Magazine Subscriptions](#)  
[Movies & TV](#)  
[Music](#)  
[Musical Instruments](#)  
[Office & School Supplies](#)  
[Patio, Lawn & Garden](#)

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

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

1.



### [Quiet: The Power of Introverts in a World That Can't Stop Talking](#)

by Susan Cain (January 24, 2012)

Average Customer Review: [★★★★★](#) (261)

In Stock

List Price: ~~\$26.00~~

Price: **\$15.60**

[110 used & new](#) from **\$14.00**

[Add to Cart](#)

[Add to Wish List](#)

☐ I own it

☐ Not interested

☒ [★★★★★](#) Rate this item

Recommended because you added **Introverts in the Church** to your Wishlist and more ([Fix this](#))

2.



### [Imagine: How Creativity Works](#)

by Jonah Lehrer (March 19, 2012)

Average Customer Review: [★★★★★](#) (107)

In Stock

List Price: ~~\$26.00~~

Price: **\$15.60**

[83 used & new](#) from **\$13.42**

[Add to Cart](#)

[Add to Wish List](#)

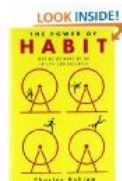
☐ I own it

☐ Not interested

☒ [★★★★★](#) Rate this item

Recommended because you added **Thinking, Fast and Slow** to your Wishlist and more ([Fix this](#))

3.



### [The Power of Habit: Why We Do What We Do in Life and Business](#)

by Charles Duhigg (February 28, 2012)

Average Customer Review: [★★★★★](#) (253)

In Stock

List Price: ~~\$28.00~~

Price: **\$15.40**

[110 used & new](#) from **\$14.91**

[Add to Cart](#)

[Add to Wish List](#)

☐ I own it

☐ Not interested

☒ [★★★★★](#) Rate this item

Recommended because you added **Thinking, Fast and Slow** to your Wishlist and more ([Fix this](#))



[Browse](#)
[Movies](#)
[Upload](#)

[Videos](#)
[Music](#)
[Movies](#)
[Shows](#)
[Live](#)
[Sports](#)
[Education](#)
[News](#)

All Categories

Recommended for You

Autos & Vehicles
Comedy
Entertainment
Film & Animation
Gaming
Howto & Style
Nonprofits & Activism
People & Blogs
Pets & Animals
Science & Technology
Travel & Events

Recommended for You »

▶ Play all

Tron Legacy lightcycle Scene
196,304 views | 1 year ago
lukasmaly7774

Daft Punk-Derezzed
330,209 views | 1 year ago
96ihatehaters

Google I/O 2011: Android + App Engine: A Developer's
48,593 views | 1 year ago
GoogleDevelopers

TRON - Daft Punk - End Of Line (Gem Siren Scene)
58,855 views | 1 year ago
triCON752

Pocoyo - 49 - Pocoyo-lympics - ENGLISH
4,306,937 views | 3 years ago
pocoyo4kids

TRON: LEGACY - "Sam Meets Castor" - Clip
309,769 views | 1 year ago
tron

The Battle of Yavin Part 1 of 2
52,388 views | 3 years ago
darkpraying

The IT Crowd - Series 1 - Episode 3: Lonely hearts
215,717 views | 3 years ago
ITCrowdChannel

Star Wars - Fire When Ready
92,388 views | 2 years ago
OneMinuteGalactica

Google I/O 2011: Android Protips: Advanced Topics
91,305 views | 1 year ago
GoogleDevelopers

Pocoyo - 67 - Duck Stuck - Season 2-15 - ENGLISH
2,958,117 views | 3 years ago
pocoyo4kids

DAFT PUNK Poses at "TRON: Legacy" World
96,418 views | 1 year ago
maximotv



# YouTube's Recommendation Algorithm Has a Dark Side

It leads users down rabbit holes








By Zeynep Tufekci on April 1, 2019 [أعرض هذا باللغة العربية](#)

These “recommended” videos play one after the other. Maybe you finished a tutorial on how to sharpen knives, but the next one may well be about why feminists are ruining manhood, how vaccinations are poisonous or why climate change is a hoax—or a nifty explainer “proving” the *Titanic* never hit an iceberg.

Zeynep Tufekci, [YouTube's Recommendation Algorithm Has a Dark Side](#), *Scientific American*, 2019


[Stories](#)
[Activity](#)

**Who to follow**

[Find friends](#)
[Browse categories](#)

United States trends · [Change](#)

[#DallasTNT](#)  Promoted

[#ConservatismIn4Words](#)

[#BiggestLiesGirlsSay](#)

[#EmailSci](#)

[Varela](#)

[And Portugal](#)

[Ronaldo](#)

[Group B](#)


[Bendtner](#)

[Danes](#)



© 2012 Twitter [About](#) [Help](#) [Terms](#) [Privacy](#) [Blog](#) [Status](#) [Apps](#) [Resources](#) [Jobs](#) [Advertisers](#) [Businesses](#) [Media](#) [Developers](#)

## Who to follow


Twitter accounts suggested for you based on who you follow and more.






**genevahenry** @genevahenry
 


Followed by cathy marshall and Richard Furuta.






**Mike Baur** @mbaur 






*I enjoy life.*  
Followed by John Stone and Stacey Vaughn.





**Bing**  @bing
 


*Bing is for doing. Tweeps: @missbeaux (^nb), @magoolovesfood (^ma), & @wangtanya (^tw)*  
Followed by Robert Scoble.  
 Promoted





**Fred Stutzman** @fstutzman
 

*Postdoc at CMU studying social networks, economics of privacy, HCI. Also developer of Freedom, Anti-Social, ClaimID.*  
Followed by Dr. Boonthum-Denecke, Eytan Adar and cathy marshall.



**Paul André** @paulesque
 

*Social Computing Postdoc at Carnegie Mellon,*

# eHarmony®

## Free to Review Your Matches

Already on Facebook?

 [Connect with Facebook](#)

First Name:

*First name only please!*

I'm a:

seeking

Zip Code:

Country:

Email:

*Note: Your email is used to log back in*

Confirm Email:

Password:

*Must be at least 5 characters*How did you  
hear about us?**Find My Matches**

## Date Smarter, Not Harder

You're not looking for lots of dates, just better ones. And that's where we come in. Take our Relationship Questionnaire to define what you are looking for, and we'll help you find the most promising matches. People whose goals, values and personality traits most complement you. **Let's get started.**

**Erika**

eHarmony member

# Collaborative Filtering

- We often seek recommendations from people we trust (family, friends, colleagues, experts)
- But who to ask for a recommendation is also based on similarity in taste
  - I trust my sister with my life, but she doesn't like the same movies I like
  - People with my tastes are likely to recommend things I will like
- CF searches a large group of people for those that like the things you like, and it combines the things they like into a ranked list of suggestions

# Web Science: Collective Intelligence & Recommender Systems (Part 3 - Recommending a Movie)

CS 432/532

Old Dominion University

*Permission has been granted to use these slides from Frank McCown, Michael L. Nelson, Alexander Nwala, Michele C. Weigle*



This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License](https://creativecommons.org/licenses/by-nc-sa/3.0/)

# Main reference this week:

Ch 2 from [Programming Collective Intelligence](#) by Toby Segaran

*(abbreviated as PCI)*

[GitHub repo](#)

# Example Data: Movie Ratings

	Lady in the Water	Snakes on a Plane	Just My Luck	Superman Returns	You, Me and Dupree	The Night Listener
Rose	2.5	3.5	3.0	3.5	2.5	3.0
Seymour	3.0	3.5	1.5	5.0	3.5	3.0
Puig		3.5	3.0	4.0	2.5	4.5
LaSalle	3.0	4.0	2.0	3.0	2.0	3.0
Matthews	3.0	4.0		5.0	3.5	3.0
Toby		4.5		4.0	1.0	

Example from Ch 2 of *PCI* (data in code on pg. 8)

Let's visualize the data on a scatterplot...

# Ratings for Snakes... and ...Dupree

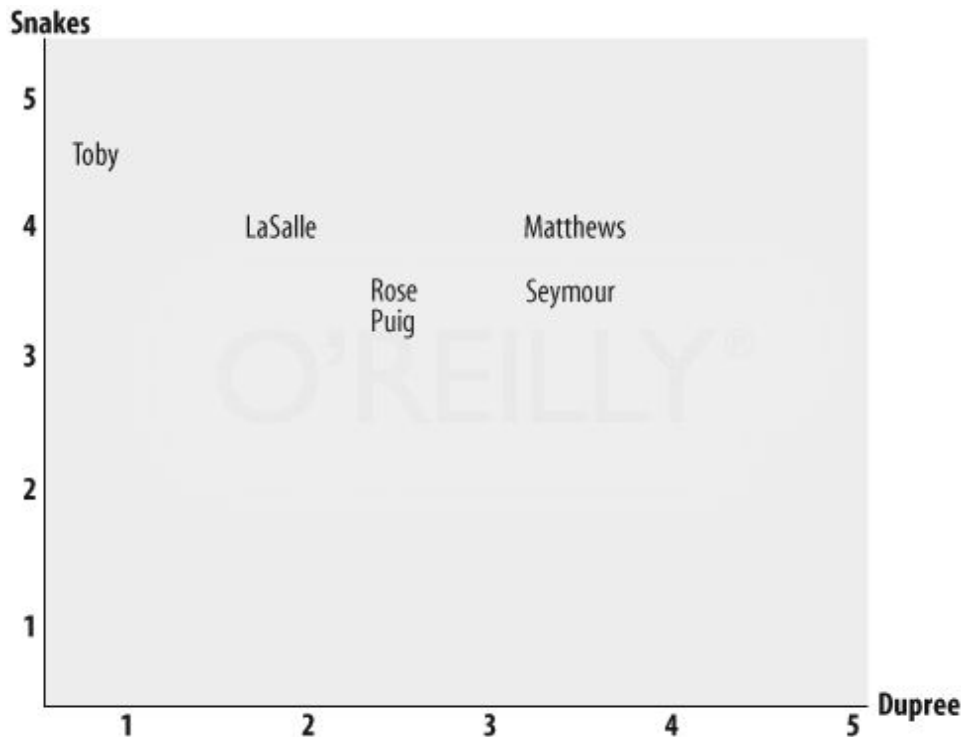
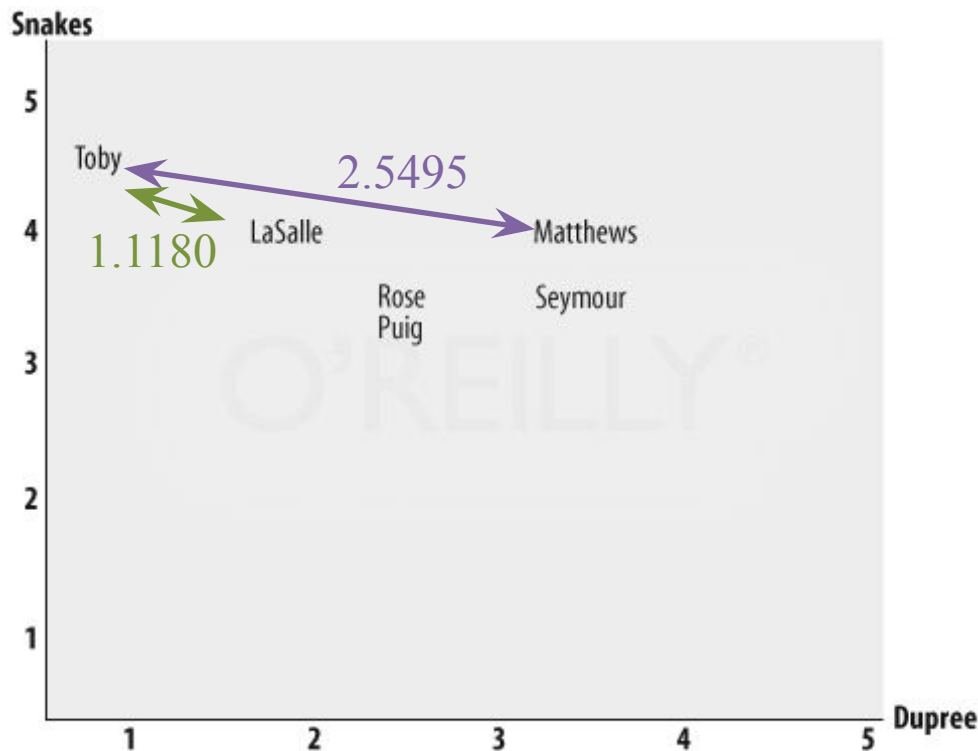


Fig 2-1, PCI, ["How to Find Similar Users with Python"](#), archived page



# Euclidean Distance Score



$$D(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

$$\begin{aligned} D(\text{Toby}, \text{Matthews}) &= \sqrt{(1 - 3.5)^2 + (4.5 - 4)^2} \\ &= 2.5495 \end{aligned}$$

$$\begin{aligned} D(\text{Toby}, \text{LaSalle}) &= \sqrt{(1 - 2)^2 + (4.5 - 4)^2} \\ &= 1.1180 \end{aligned}$$

PCI, pgs. 8-15, ["How to Find Similar Users with Python"](#), archived page

# Modified Euclidean Distance Score

- Most similarity metrics use range [0,1]  
0 = no similarity  
1 = exactly the same
- Modify Euclidean Distance (to measure similarity)

$$D(x, y) = 1 / (1 + \sqrt{\sum_{i=1}^n (x_i - y_i)^2})$$

# Modified Euclidean Distance Score

	Lady in the Water	Snakes on a Plane	Just My Luck	Superman Returns	You, Me and Dupree	The Night Listener
LaSalle	3.0	4.0	2.0	3.0	2.0	3.0
Matthews	3.0	4.0		5.0	3.5	3.0
Toby		4.5		4.0	1.0	

Distance(Toby, LaSalle) - original distance value was 1.1180

$$1/(1 + \sqrt{(4.5 - 4)^2 + (4 - 3)^2 + (1 - 2)^2}) = 0.4$$

Distance(Toby, Matthews) - original distance value was 2.5495

$$1/(1 + \sqrt{(4.5 - 4)^2 + (4 - 5)^2 + (1 - 3.5)^2}) = 0.2675$$

# Pearson Correlation Coefficient

- Problem with Euclidean distance

Ratings for A: 5, 4, 3      Ratings for B: 3, 2, 1

*although A and B seem to share roughly same opinions, distance between ratings are significant*

- Pearson's  $r$  corrects for "grade inflation"
- Yields values between 1 and -1
  - 1 = perfect correlation
  - 0 = no correlation
  - 1 = inverse correlation

# Calculating Pearson's $r$

$x$  = User 1's ratings of all items rated by both users

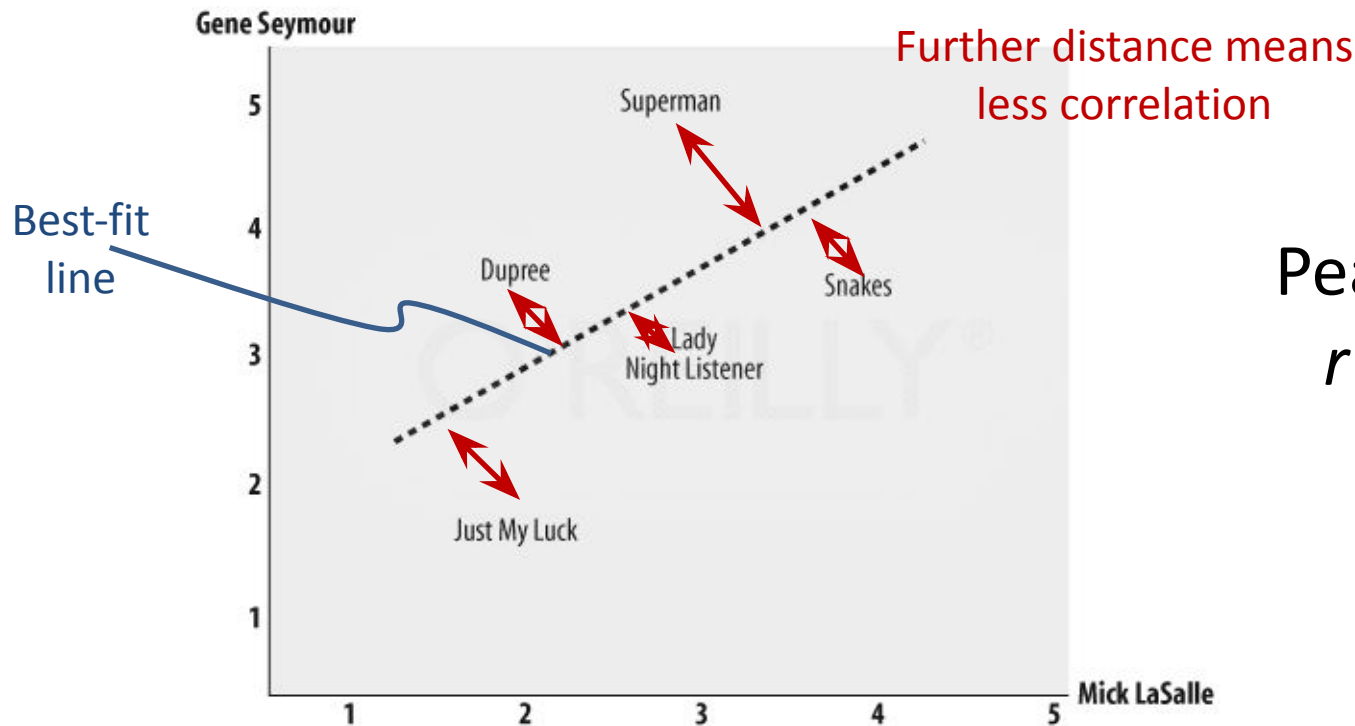
$y$  = User 2's ratings of all items rated by both users

$\bar{x}$  = mean of all User 1's ratings

$\bar{y}$  = mean of all User 2's ratings

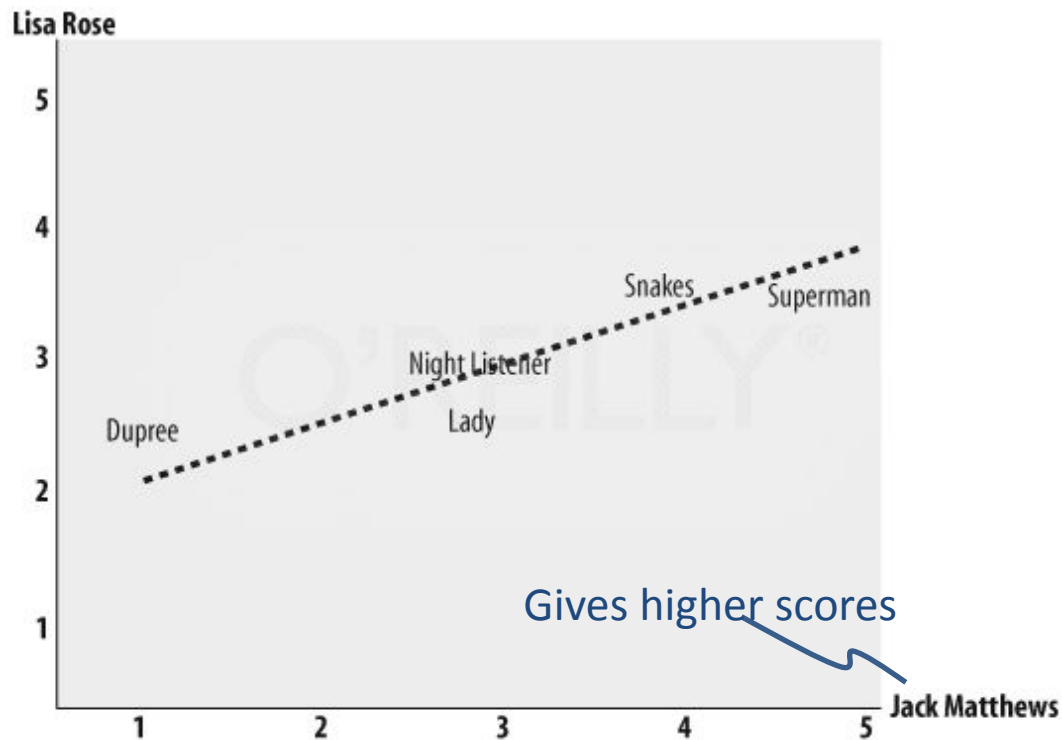
$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

# Comparing Two Critics' Ratings



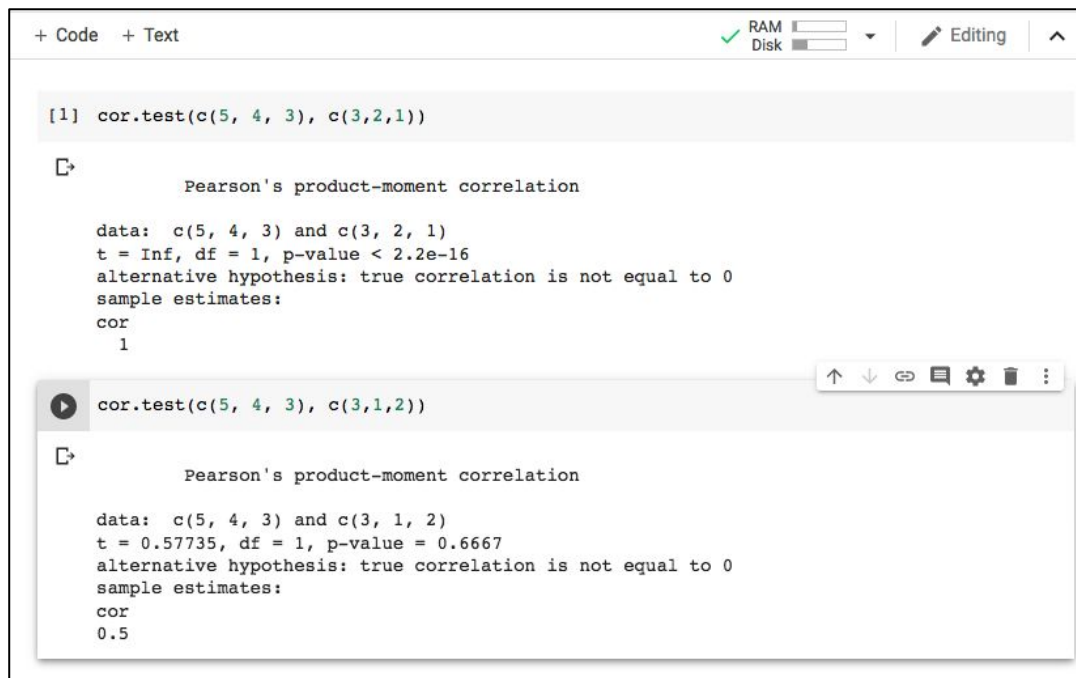
Pearson's  
 $r = 0.4$

# Comparing Two Critics' Ratings



Pearson's  
 $r = 0.75$

# Pearson's $r$ in R



The screenshot shows a Google Colab notebook interface. At the top, there are tabs for '+ Code' and '+ Text', and a status bar indicating 'RAM' and 'Disk' usage, along with an 'Editing' mode icon. The first code cell contains the R command `cor.test(c(5, 4, 3), c(3,2,1))`. Its output displays 'Pearson's product-moment correlation' with data points `c(5, 4, 3)` and `c(3, 2, 1)`, a `t = Inf` value, `df = 1`, a `p-value < 2.2e-16`, and a sample estimate of `cor = 1`. The second code cell contains the R command `cor.test(c(5, 4, 3), c(3,1,2))`. Its output displays 'Pearson's product-moment correlation' with data points `c(5, 4, 3)` and `c(3, 1, 2)`, a `t = 0.57735` value, `df = 1`, a `p-value = 0.6667`, and a sample estimate of `cor = 0.5`. A toolbar with icons for undo, redo, link, comment, settings, and delete is visible between the two code cells.

```
[1] cor.test(c(5, 4, 3), c(3,2,1))
```

Pearson's product-moment correlation

data: c(5, 4, 3) and c(3, 2, 1)  
t = Inf, df = 1, p-value < 2.2e-16  
alternative hypothesis: true correlation is not equal to 0  
sample estimates:  
cor  
1

```
cor.test(c(5, 4, 3), c(3,1,2))
```

Pearson's product-moment correlation

data: c(5, 4, 3) and c(3, 1, 2)  
t = 0.57735, df = 1, p-value = 0.6667  
alternative hypothesis: true correlation is not equal to 0  
sample estimates:  
cor  
0.5

[Create new R Google Colab notebook](#)

Python: [scipy.stats.pearsonr](#) — [SciPy v1.5.3 Reference Guide](#), and pg. 13 in *PCI*



# Other Similarity Measures

- [Cosine similarity](#) (Wikipedia)
- [Jaccard coefficient](#) (Wikipedia)
- [Manhattan \(taxicab\) distance](#) (Wikipedia)
- Others...

# Moving On...

	Lady in the Water	Snakes on a Plane	Just My Luck	Superman Returns	You, Me and Dupree	The Night Listener
Rose	2.5	3.5	3.0	3.5	2.5	3.0
Seymour	3.0	3.5	1.5	5.0	3.5	3.0
Puig		3.5	3.0	4.0	2.5	4.5
LaSalle	3.0	4.0	2.0	3.0	2.0	3.0
Matthews	3.0	4.0		5.0	3.5	3.0
Toby		4.5		4.0	1.0	

Should Toby see these movies?

# Who Should We Ask?

- To find recommendations for movies we have not seen, we could...
  1. Find a *single* critic whose taste best matches ours  
*What if they haven't rated a movie we are interested in?*
  2. Get *all* critics' input but give critics with *similar* tastes more impact on decision
- For option 2, use similarity metric to compare all ratings with ours, and compute average rating based on weighted similarity

# Should Toby See "Lady in the Water"?

All Toby's ratings compared to other critics'

Higher values influence Weighted Avg more

Not included in Total since no rating for Lady in the Water

	Similarity (Pearson's $r$ )	Lady in the Water	Weighted Score
Rose	0.99	2.5	2.48
Seymour	0.38	3.0	1.14
Puig			
LaSalle	0.92	3.0	2.77
Matthews	0.66	3.0	1.99
Total	2.95		8.38

$0.99 \times 2.5$

Weighted Avg = Weighted Total / Similarity Total =  $8.38 / 2.95 = 2.83$

Toby is likely to rate this movie 2.83

# Product Recommendations

- Previous example used distance metric to find critics with similar taste
- What if we want to find movies that are similar to some given movie?
- Solution: Use same method but swap rows and columns

# Rows and Columns Swapped

	Rose	Seymour	Puig	LaSalle	Matthews	Toby
Lady in the Water	2.5	3.0		3.0	3.0	
Snakes on a Plane	3.5	3.5	3.5	3.5	4.0	4.0
Just My Luck	3.0	1.5	3.0	2.0		
Superman Returns	3.5	5.0	4.0	3.0	5.0	4.0
You, Me and Dupree	2.5	3.5	2.5	2.0	3.5	1.0
The Night Listener	3.0	3.0	4.5	3.0	3.0	

Find movies like *Superman Returns* by comparing its row with all other rows...

# Movies Similar to *Superman Returns*

	Similarity
You, Me and Dupree	0.657
Lady in the Water	0.487
Snakes on a Plane	0.111
The Night Listener	-0.179
Just My Luck	-0.422

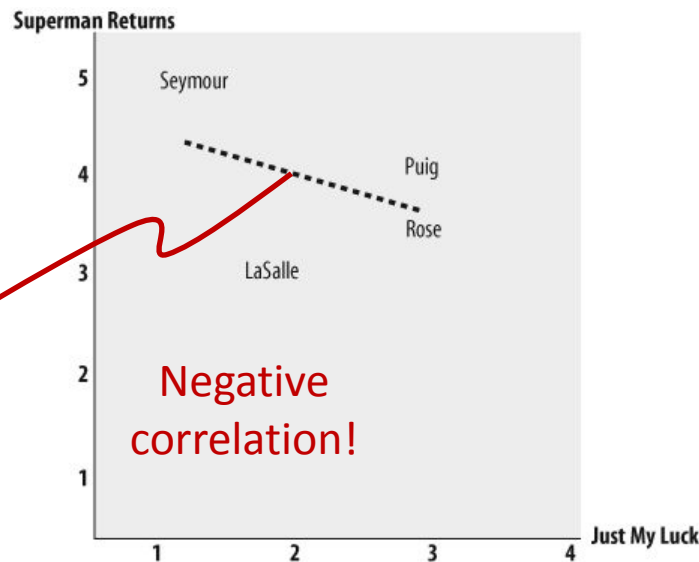


Fig 2-5 from PCI

If you like *Superman Returns*, you probably won't like *Just My Luck* (and vice versa)!

# Web Science: Collective Intelligence & Recommender Systems

(Part 4 - Challenges for Collaborative Filtering)

CS 432/532

Old Dominion University

*Permission has been granted to use these slides from Frank McCown, Michael L. Nelson, Alexander Nwala, Michele C. Weigle*



This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License](https://creativecommons.org/licenses/by-nc-sa/3.0/)



# Challenges for Collaborative Filtering

- Ratings data is often sparse
  - Large data sets required
  - Cold start problem: new users must first rate a number of items before CF can work
- Comparing all items against all users is not efficient/scalable
  - Compare against a sample
  - Use clustering algorithm to locate best neighbors

Clustering algorithms are covered in Ch 3 of [\*Programming Collective Intelligence\*](#)

# Challenges for Collaborative Filtering

- Susceptible to cheating (Shilling attacks)
  - Rate your own products higher and your competitors lower
  - Rate like everyone else except a select few items you want recommended to others



Lam and Riedl, "[Shilling Recommender Systems for Fun and Profit](#)", WWW 2004

Img source: <http://www.somethingawful.com/d/photoshop-phriday/cheating-in-sports.php>

# Say it isn't true!

## Hidden Industry Dupes Social Media Users

Paying people to influence discussions in social media is big business in China and the U.S.

By Tom Simonite on December 12, 2011

A trawl of Chinese crowdsourcing websites – where people can earn a few pennies for small jobs such as labeling images – has uncovered a multimillion-dollar industry that pays hundreds of thousands of people to distort interactions in social networks and to post spam.

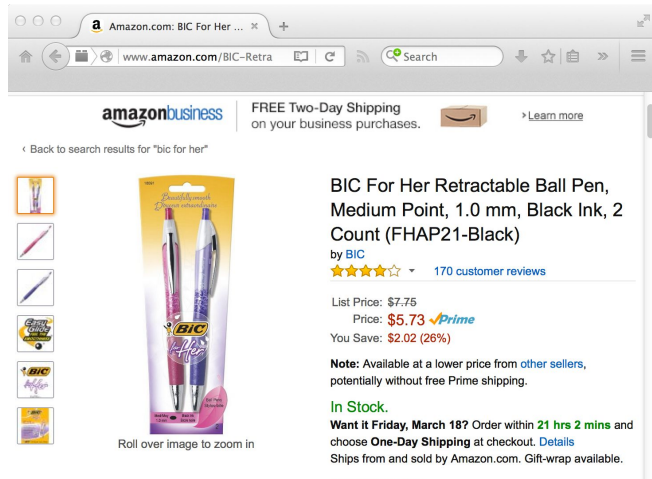


[Hidden Industry Dupes Social Media Users](#)

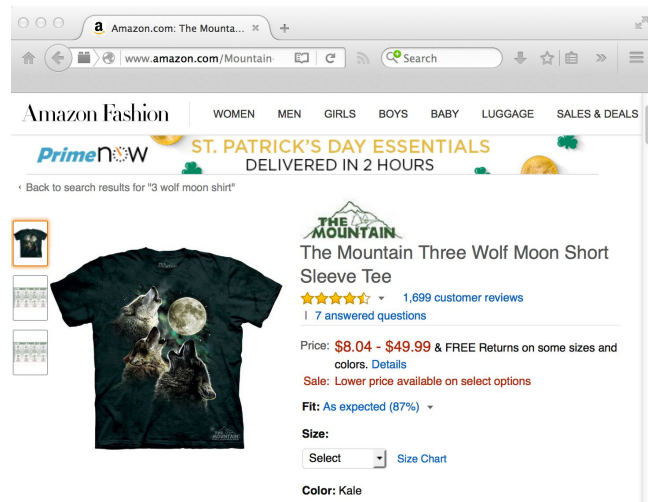
[Astroturfing](#) (Wikipedia)

[OSoMe](#) (Observatory on Social Media), Indiana University

# Some reviews are sarcastic...



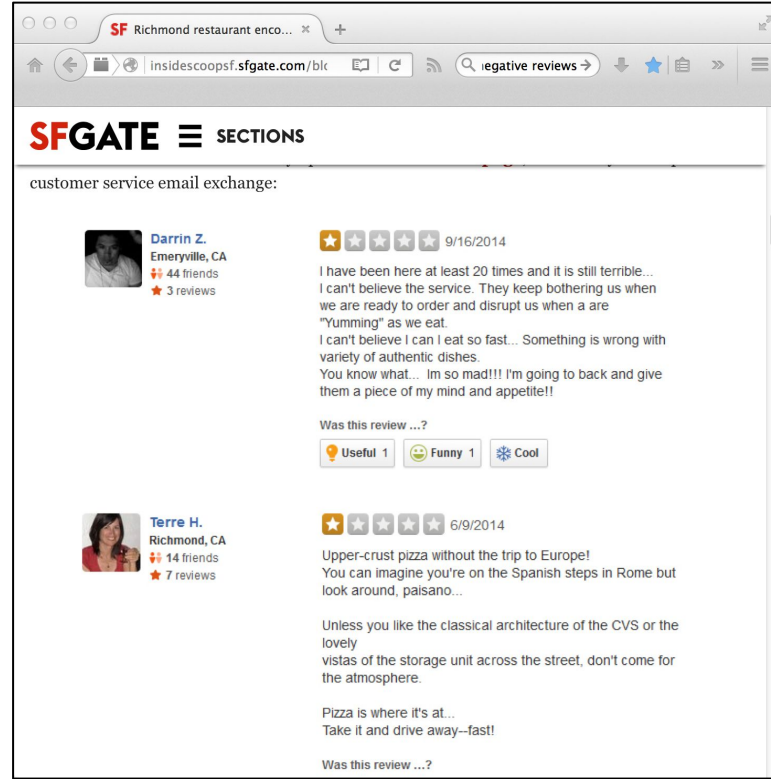
5.0 out of 5 stars Received As Wedding Gift- LOVE THEM!  
By agilerton on August 28, 2012  
Style Name: BlueSize: 1 - Pack  
I received a set of these pens for my wedding. Unfortunately I can't write or type (my husband is actually typing for me.) However, I do so love the pretty colors. I've been learning how to write with them. ...



5.0 out of 5 starsGreat compliment for my skin art  
By overlook1977 on May 19, 2009  
Size: MediumColor: Dark Green  
Unfortunately I already had this exact picture tattooed on my chest, but this shirt is very useful in colder weather.

[The Mountain Men's Three Wolf Moon Short Sleeve Tee at Amazon.com](https://www.amazon.com/dp/B000APR000)

# Some are protests.



[Richmond restaurant encourages bad Yelp reviews](#)  
[Botto Italian Bistro](#) (yelp)

# Netflix Prize

- Netflix Prize (Oct 2006): \$1M for beating Netflix's collaborative filtering algorithm by 10%
- Dataset: 100,480,507 ratings that 480,189 anonymized users gave to 17,770 movies
- Started with 50,051 contestants
- Only two teams had winning solutions by contest closing (July 2009)



Img source: [http://www.wired.com/images\\_blogs/business/2009/09/p1010915.jpg](http://www.wired.com/images_blogs/business/2009/09/p1010915.jpg)

# How Did They Do It?

- Insights that gave winning algorithms an edge:
  - People who rate a large number of movies at once are usually rating movies they saw a long time ago
  - People use different criteria when rating movies they saw a long time ago vs. recently seen
  - Some movies get better over time, others get worse
  - People rate movies differently depending on which day of the week it is
- Combined hundreds of other algorithms which were precisely weighed and tuned

[How the Netflix Prize Was Won](#) by Buskirk, Wired 2009

# Netflix Prize 2

- Second prize announced in 2009
- FTC and lawsuits against Netflix about privacy
- I thought the data was **anonymized?**



Img source: [http://laist.com/2008/02/10/photo\\_essay\\_ano.php](http://laist.com/2008/02/10/photo_essay_ano.php)



# Anonymizing Data Is Hard To Do...

- In 2007 Narayanan and Shmatikov (Univ of Texas) were able to re-identify a number of the anonymous users in Netflix dataset by correlating movie ratings with IMDb ratings ([full paper](#))
- This has happened before...
  - In 2006 users in anonymized AOL search data were re-identified
  - In 2000 Latanya Sweeney showed 87% of all Americans could be uniquely identified with only zip code, birthdate, and gender

See [Why Anonymous Data Sometimes Isn't](#) by Bruce Schneier (2007)

Netflix gave up...  
no more prizes!



# Objectives

- Describe and define the four components of collective intelligence.
- Explain how collaborative filtering is related to recommender systems.
- Differentiate between Euclidean distance and the Pearson correlation coefficient.
- List three challenges for collaborative filtering.