# Social Data Mining

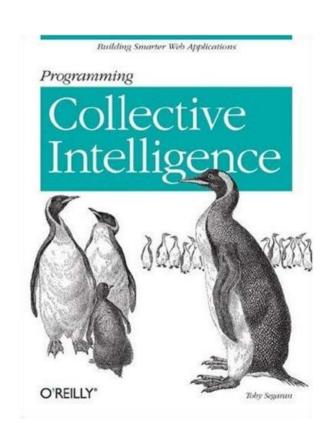
**Toby Segaran** 

### About Me







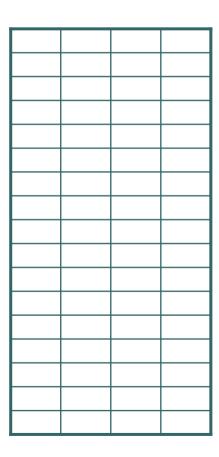


http://kiwitobes.com

# • • What is data mining?

- Implicit
- Unknown
- Useful

#### What is data?



.1

COJ: COUNCIL 2004-06-...?
COJ: MAYORAL COMMITTEE 2004-06-10
COJ: FINANCE: STRATEGY & ECONOMIC DEVELOPMENT COMMITTEE 2004-06-07

#### (ECONOMIC DEVELOPMENT UNIT)

\_\_\_\_

... URBAN RENEWAL TAX INCENTIVE SUBMISSION TO NATIONAL TREASURY

1. STRATEGIC THRUST

Economic Growth and Development, Housing Delivery, Service Delivery Excellence, Inner City Regeneration and Good Governance

2. OBJECTIVE

To inform committees on the City of Johannesburg's (CoJ) application to National Treasury for the designation of the inner City as Johannesburg's Urban Development Zonc (UDZ) in line with the National Treasury's urban renewal tax incentive.

3. BACKGROUND

Legislative proposats dealing with the Urban Renewal Tax Incentive were signed into law in December 2003, i.e. "Revenue Laws Amendment Act," 2003: Insertion of section 13 quat in Act 58 of 1996 (Incente Tax Act). "Deductions in respect of erection or improvement of buildings in urban development zones".

The proposed tax incentive applies to the construction and refurbishment of both commercial and residential buildings in designated decayed or under vullised inner city areas and will be available to owners or lessors of the properties. The incentive will take the form of a special depreciation allowance for the refurbishment of existing buildings (taxpayers will receive a 20% straight-inner depreciation allowance over a 5-year period) and/or the construction of new buildings (taxpayers will receive a 17-year write-off period with a 20% write-off in the first year, and a 5% write-off thereafter).

4. JOHANNESBURG'S UDZ: THE INNER CITY

In terms of giving effect to the legislation, the CoJ has prepared an application for submission to Notional Treasury (see Annexure A) by the end of June 2004. Key elements contained in the application are outlined rext.

4.1 MOTIVATION FOR EXTENT OF UDZ AREA BEING APPLIED FOR BY COJ

The Inner City area is nearly 1,800 nectares in extent. The maximum permissible area oillowable in forms of the legislation, related to Johannsehung's population, is 690 hectares. However, an application for a larger area can be made. A detailed motivation to the Ministor of Finance for the declaration of the entire Inner City as Johannesburg's UDZ has been provided. Section A of the application describes the Inner City as a system of four main zones, each with different characteristics, i.e.:



#### Data-mining traditional uses











# • • Why it's important now



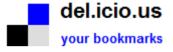
### Why it's important now































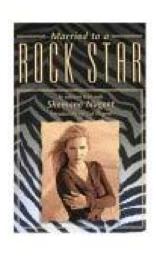




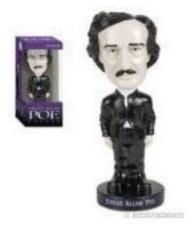
#### Why it's important now













All products are actually sold on Amazon

### Why it's important now

#### **Facebook**



Google

Sponsored Links

#### private tours of Berlin

English speaker will show you the sights of **Berlin** (or in Spanish) www.afriendinberlin.de

#### **Berlin Tours**

Learn More about Destinations All Around Europe - Find Great Info! www.BudgetTravel.com

#### **Berlin Tours**

Find and book 25+ things to do in **Berlin** on Viator. www.viator.com/berlin

#### Berlin Sightseeing Tours

Tour Berlin & Surrounding Germany. Book Online & Save. 800-208-4421 www.Berlin-Tours.net

30 colors

# • • Why it's important now







# • • | For Social Insight



Home Prices



Blogs and News



Movie Data



**Fashion** 



Hotties

#### Blogs...





All Stories Web Broadband

#### NewTeeVee Liv & Speakers



Written by Om Malik Sunday, October 14, 2007

For past few months, the entire Wan has been busy pulling tog conference, which is slotted to Email Me Francisco. We posted the sche Friday.



SETH'S WEB PAGES

The Dip Blog

Books by Seth Godin

Seth's Main Blog

Seth's Squidoo Lens

All Marketers Are Liars Blog SethGodin.com: Official Site

RSS FEEDS

#### Marketing Fear

Marketing with fear is a powerful tool. Fear is a universal emotion, it's viral and people will go to great lengths to make it go away.

Some items can't be marketed without fear. Seat belts, for example. They're not convenient, good tasting, fun to use or profitable. Fear works great in this

An essential question to ask, though, is who benefits? In the case of seat belts, the use of fear directly benefits the prospect, because using seat belts not only

decreases fe you're selling \$30 warrant investment.









news organi benefit, Cra sell some du



gothamist\*

Search Gothamist GO

tips@gothamist.com

ADVERTISEMENT





WHIT

Features: Blog | Newsmap | Contribute |

OCTOBER 14, 2007

#### **Times Weddings Highlig**



### The Technorati Top 100

Most Linked To

Most Favorited

#### 1. Engadget

By WeblogsInc · 54 minutes ago
http://www.engadget.com
Authority: 31,457

#### 2. Boing Boing: A Directory of Wonderful Things

Boing Boing is a weblog of cultural curiosities and interesting technologies. It's the most popular blog in the world, as ranked by Technorati.com, and won the Lifetime Achievement and Best Group Blog awards at the 2006 Bloggies ceremony.

By Mark Frauenfelder · 23 minutes ago

http://www.boingboing.net

Authority: 23,507

#### 3. Gizmodo, the Gadget Guide

39 minutes ago

Mttp://gizmodo.com

Authority: 22,172

#### 4. Techcrunch



### • Getting the content



#### Thu The Six Degrees Hypothesis Experienced

It is when you travel that you most powerfully experience the reality of Stanley Milgr hypothesis, the idea that we're only six degrees away from anyone else.

You'd think that with billions of people in the world, the chance of you running into : if only indirectly) would be tiny, yet in my travels, it's happened to me repeatedly.

I'm on holiday in Sicily. I'm walking down the streets of Ortygia, the old town of Sirac say, "My god, it's Tim O'Reilly." It's Kevin Altis, whom I've never met, but who helps track at Oscon. He too is on holiday.

That night, at dinner, at a fabulous restaurant called Don Camillo, where we've gor in salt, we see a large group at the next table, all speaking English, but with a coup Italian. We ask for a bit of help with the menu. We discover that this is a group of fc food industry executives on a tour with Boulder-based Culinary Adventures.

The

Six

Degrees

**Hypothesis** 

Experienced

lt

ls

When

You

Travel







### Building a Word Matrix

The

Six

Degrees

Hypothesis

Experienced

lt

Is

When

You

Travel

Six

Degrees

Hypothesis

Experienced

Travel



Six	3
Degrees	3
Hypothesis	1
Experienced	5
Travel	6

# • • The Word Matrix

	"china"	"kids"	"music"	"travel"	"yahoo"
Gothamist	0	3	3	3	0
GigaOM	6	0	1	4	2
QuickOnlineTips	0	2	2	0	12
O'Reilly Radar	1	0	3	6	4

# • • Determining distance

	"china"	"kids"	"music"	"yahoo"
Gothamist	0	3	3	0
GigaOM	6	0	1	2
Quick Online Tips	0	2	2	12

#### Euclidean "as the crow flies"

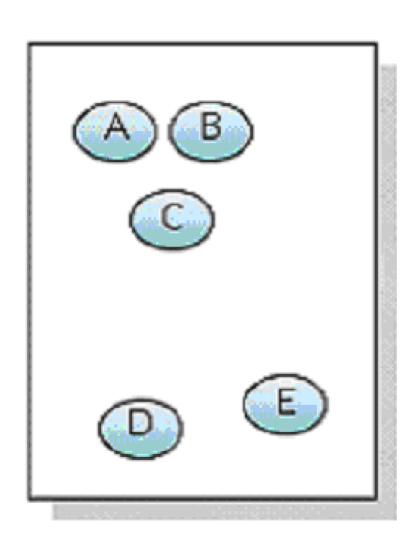
$$\sqrt{(6-0)^2+(0-2)^2+(1-2)^2+(2-12)^2}$$

= 12 (approx)

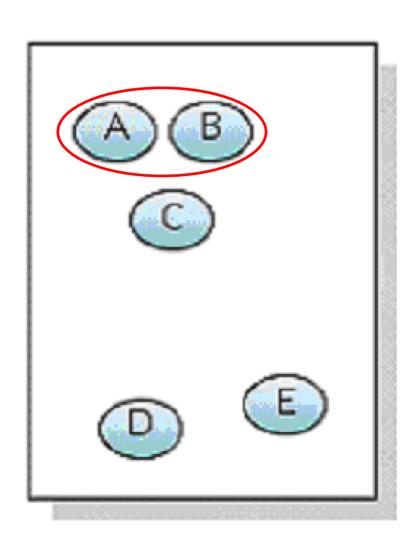
# • • Hierarchical Clustering

- Find the two closest item
- Combine them into a single item
- Repeat...

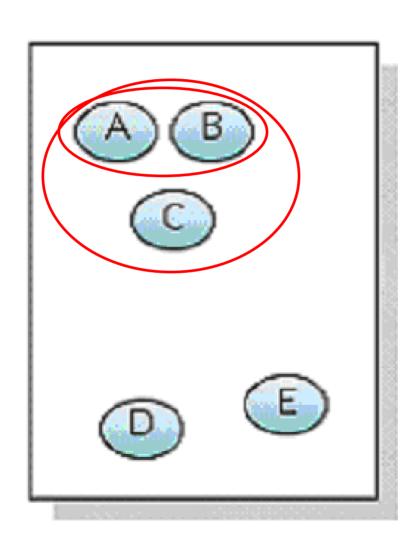
### Hierarchical Algorithm



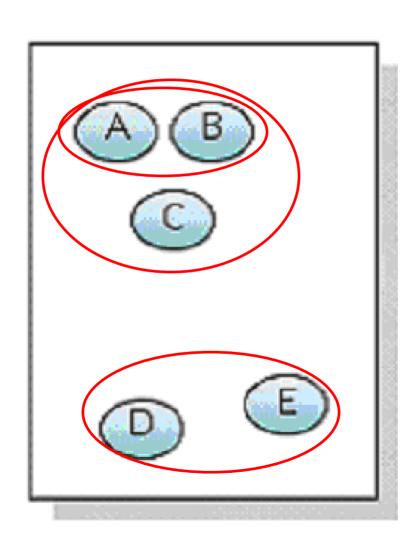
### Hierarchical Algorithm



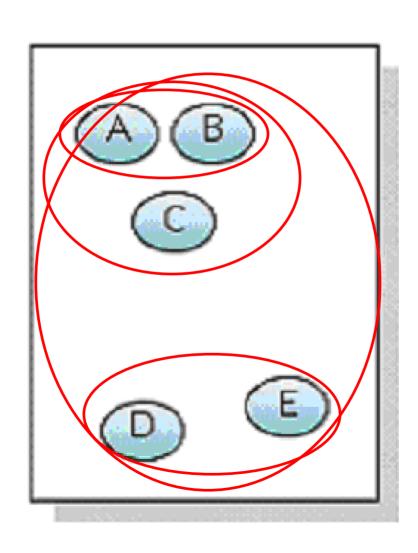
### Hierarchical Algorithm



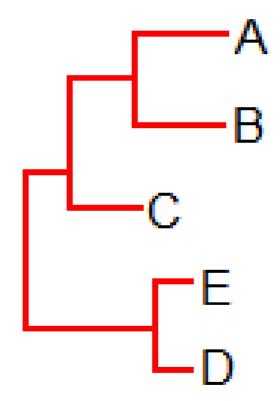
### • • Hierarchical Algorithm



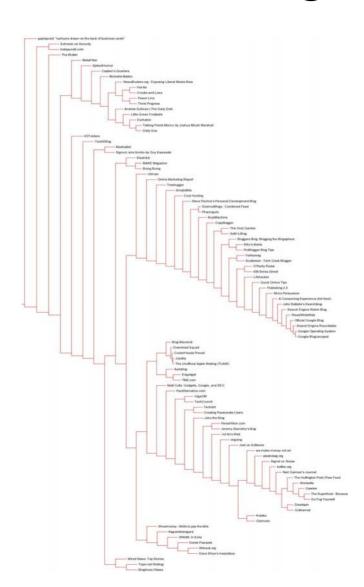
# • • Hierarchical Algorithm



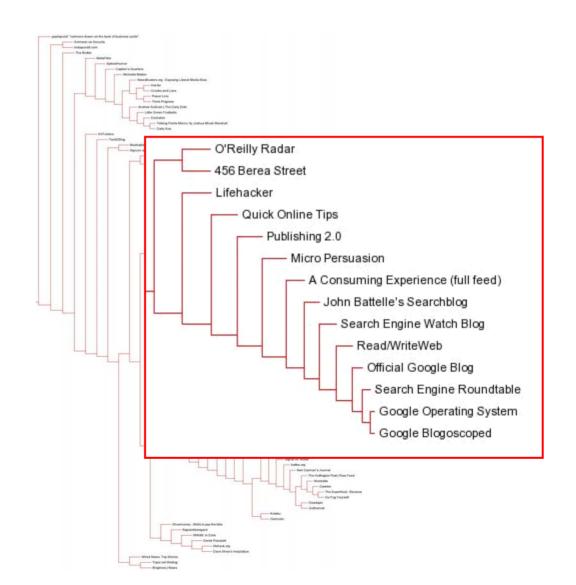
# • • Dendrogram



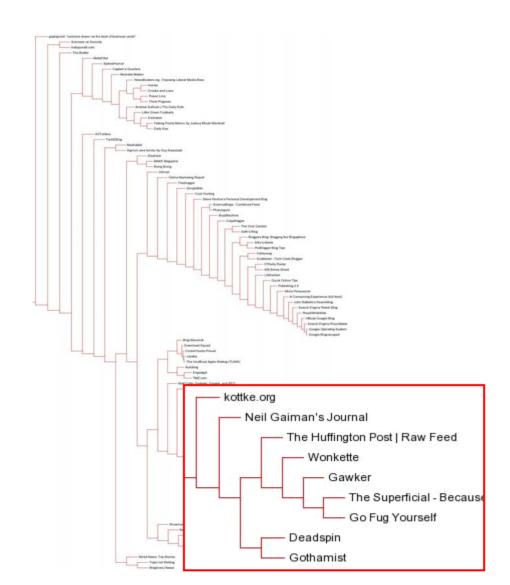
### • • Hierarchical Blog Clusters



### Hierarchical Blog Clusters



## • • Hierarchical Blog Clusters

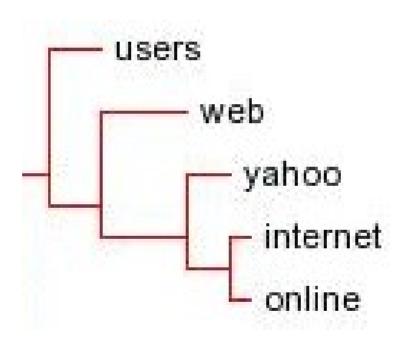


# Rotating the Matrix

#### Words in a blog -> blogs containing each word

	Gothamist	GigaOM	Quick Onl
china	0	6	0
kids	3	0	2
music	3	1	2
Yahoo	0	2	12

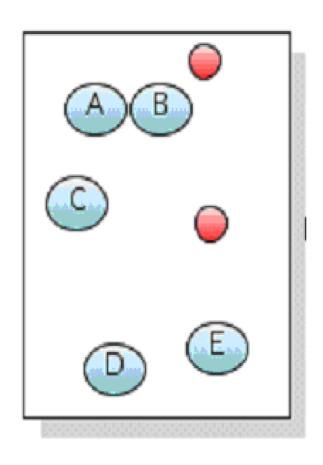
### • • Hierarchical Word Clusters



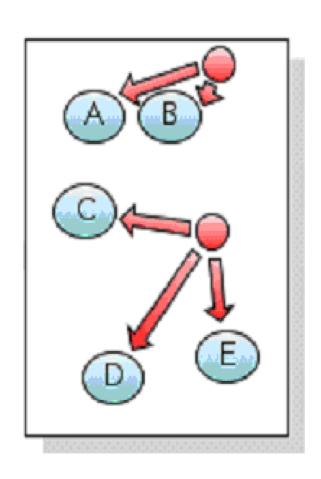
# K-Means Clustering

- Divides data into distinct clusters
- User determines how many
- Algorithm
  - Start with arbitrary centroids
  - Assign points to centroids
  - Move the centroids
  - Repeat

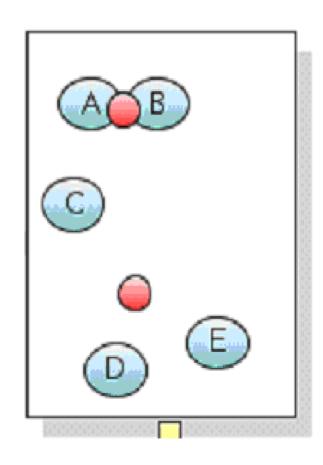
# • • K-Means Algorithm



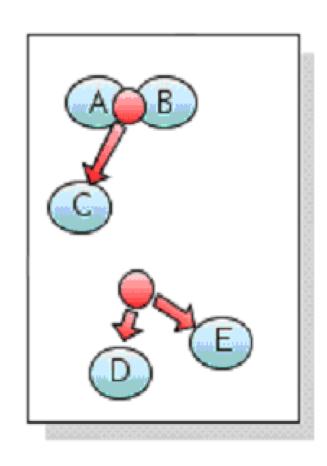
### K-Means Algorithm



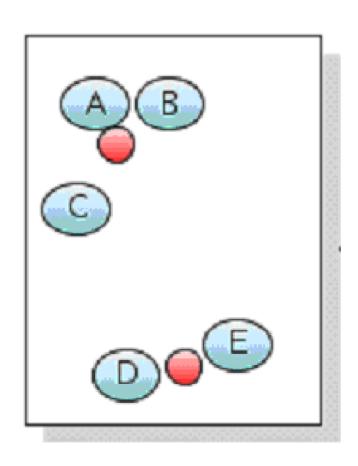
# • • K-Means Algorithm



### K-Means Algorithm



### K-Means Algorithm



### • • K-Means Results

1

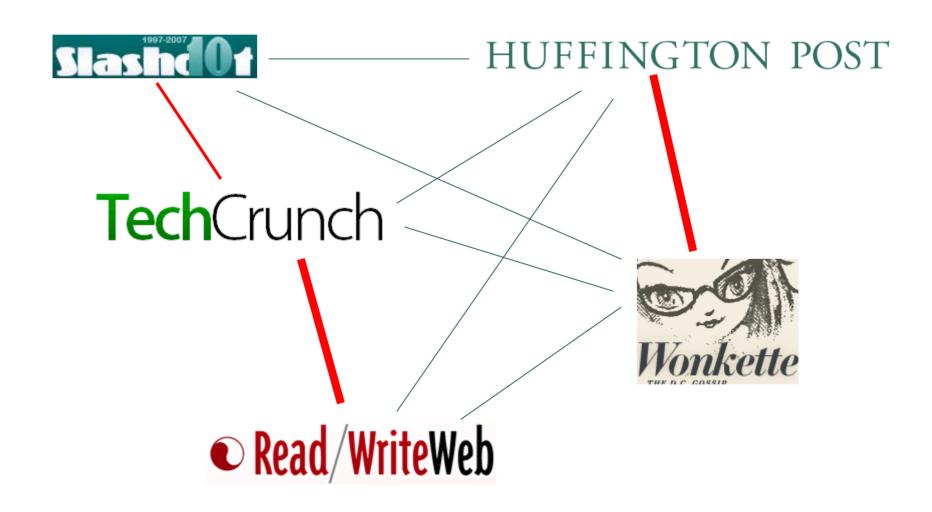
The Viral Garden
Copyblogger
Creating Passionate Users
Oilman
ProBlogger Blog Tips
Seth's Blog

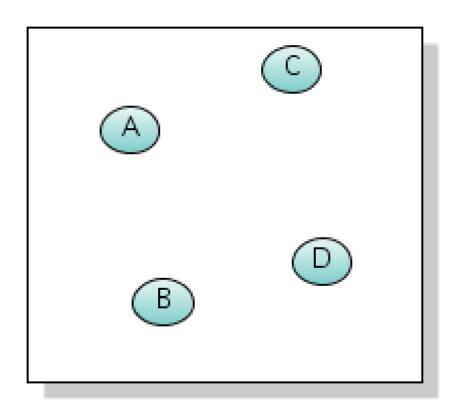
2

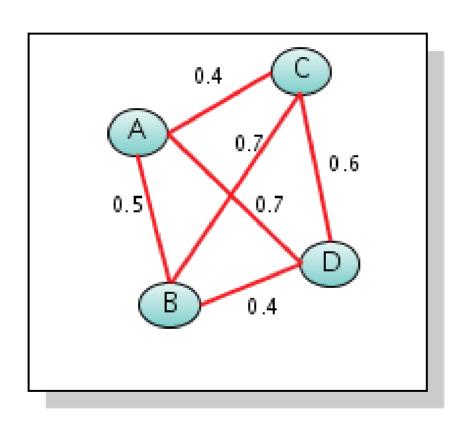
Wonkette
Gawker
Gothamist
Huffington Post

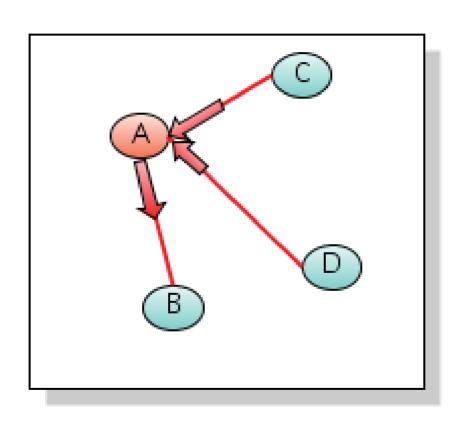
### • • 2D Visualizations

- Instead of Clusters, a 2D Map
- Goals
  - Preserve distances as much as possible
  - Draw in two dimensions
- Dimension Reduction
  - Principal Components Analysis
  - Multidimensional Scaling

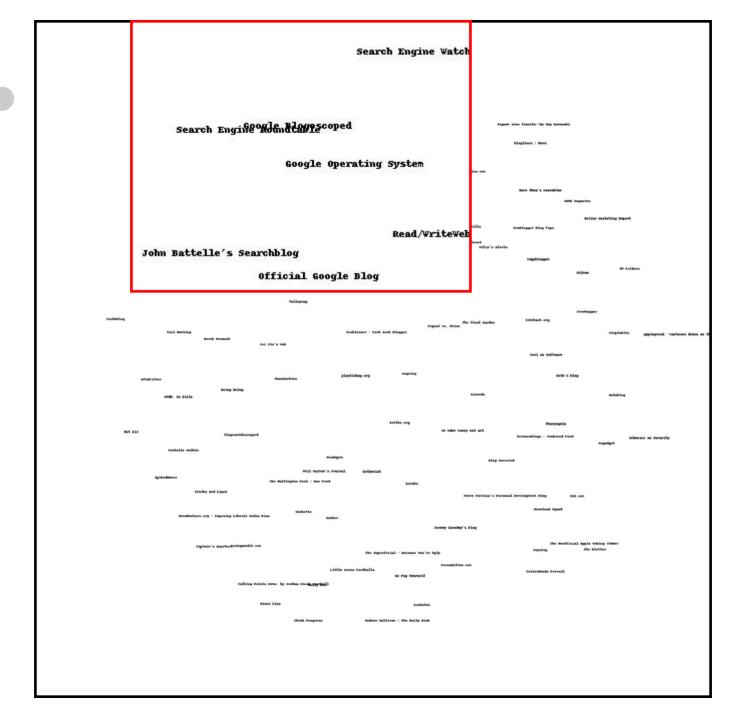








0'Reilly Radar official teeple Bloy Jeel on Seffware web: In Exile Not Air Talking Points Nemo: by Joshua Micaharnyshall



The Huffington Post | Raw Crooks and Liars Wonkette ewsBusters.org - Exposing Liberal Media Bias

Techdirt Topix.net Weblog Slashdot Wired News: Top Stories **TechEBlog** 

# • • Zillow

Home Facts			
Public Facts	;	Owner's Facts	
Residence:	Multi family	The owner has not	
Bedrooms:	5	edited home facts or created an estimate.	
Bathrooms:	3.5	_	
Sq ft:	2,474	Are you the owner?  Edit home facts	
Lot size:	2,819 sq ft / 0.06 acres	Learn more	
Year built:	1902	Create an estimate	
Year updated:		After you're done, your	
# Stories:	3	edited home facts will appear here. You can also	
# Units:	3	make your estimate public	
Total rooms:	11	or keep it private.	
Zestimate:	\$557,447		
<b>⊕</b> Show all home facts			

### • • The Zillow API

- Allows querying by address
- Returns information about the property
  - Bedrooms
  - Bathrooms
  - Zip Code
  - Price Estimate
  - Last Sale Price

## • • A home price dataset

House	Zip	Bathrooms	Bedrooms	Built	Туре	Price
А	02138	1.5	2	1847	Single	505296
В	02139	3.5	9	1916	Triplex	776378
С	02140	3.5	4	1894	Duplex	595027
D	02139	2.5	4	1854	Duplex	552213
E	02138	3.5	5	1909	Duplex	947528
F	02138	3.5	4	1930	Single	2107871
etc						

# • • What can we learn?

- A made-up houses price
- o How important is Zip Code?
- What are the important attributes?

• Can we do better than averages?

# Introducing Regression Trees

Α	В	Value	
10	Circle	20	
11	Square	22	
22	Square	8	A < 12 ?
18	Circle	6	7 12 1
		22 Flowered by yE	B = Circle ?  No  Yes  No  Yes  No  Yes  A  B = Circle ?  B = Circle ?  A  B = Circle ?  B = Circle ?  B = Circle ?

# Introducing Regression Trees

Α	В	Value	
10	Circle	20	
11	Square	22	
22	Square	8	A < 12 ?
18	Circle	6	A 12 !
		22 Powered by yEi	B = Circle ?  No Yes  avg = 7  B = Circle ?  No Yes  8 6

- Standard deviation is the "spread" of results
- Try all possible divisions
- Choose the division that decreases deviation the most

А	В	Value
10	Circle	20
11	Square	22
22	Square	8
18	Circle	6

#### **Initially**

Average = 14

Standard Deviation = 8.2

- Standard deviation is the "spread" of results
- Try all possible divisions
- Choose the division that decreases deviation the most

А	В	Value
10	Circle	20
11	Square	22
22	Square	8
18	Circle	6

#### B = Circle

Average = 13

Standard Deviation = 9.9

Average = 15

Standard Deviation = 9.9

- Standard deviation is the "spread" of results
- Try all possible divisions
- Choose the division that decreases deviation the most

А	В	Value
10	Circle	20
11	Square	22
22	Square	8
18	Circle	6

Average 
$$= 8$$

Standard Deviation = 0

$$A <= 20$$

Average 
$$= 16$$

Standard Deviation = 8.7

- Standard deviation is the "spread" of results
- Try all possible divisions
- Choose the division that decreases deviation the most

А	В	Value
10	Circle	20
11	Square	22
22	Square	8
18	Circle	6

Average = 7

Standard Deviation = 1.4

$$A <= 11$$

Average = 21

Standard Deviation = 1.4

# • • CART Algoritm

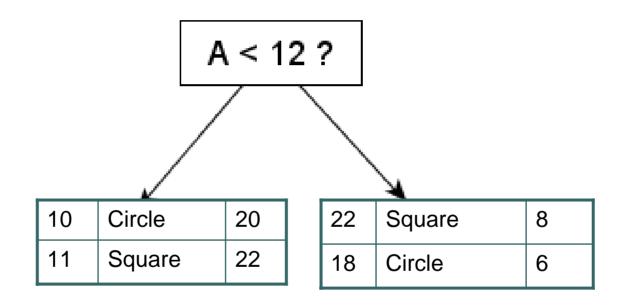
А	В	Value
10	Circle	20
11	Square	22
22	Square	8
18	Circle	6

# • • CART Algoritm

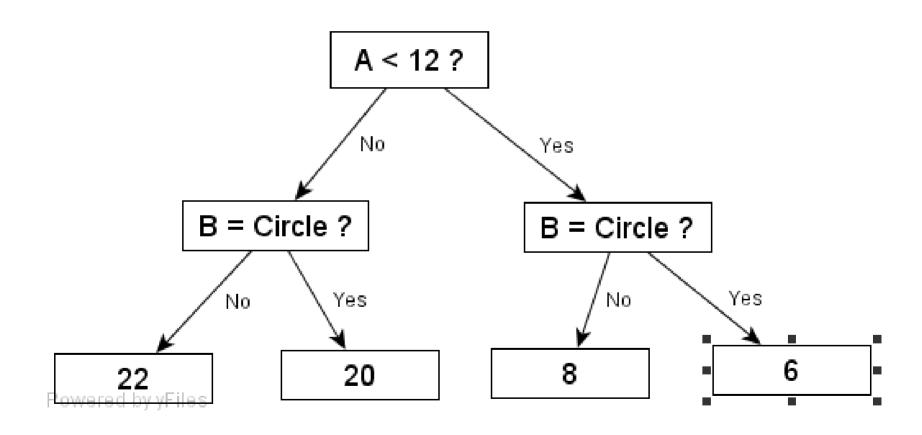
А	В	Value
10	Circle	20
11	Square	22
22	Square	8
18	Circle	6

A < 12 ?

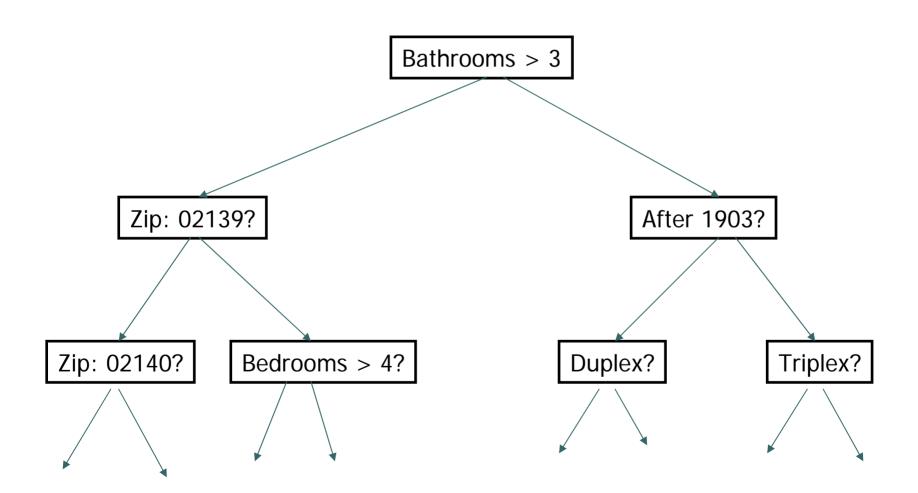
### • • • CART Algoritm



### • • CART Algoritm

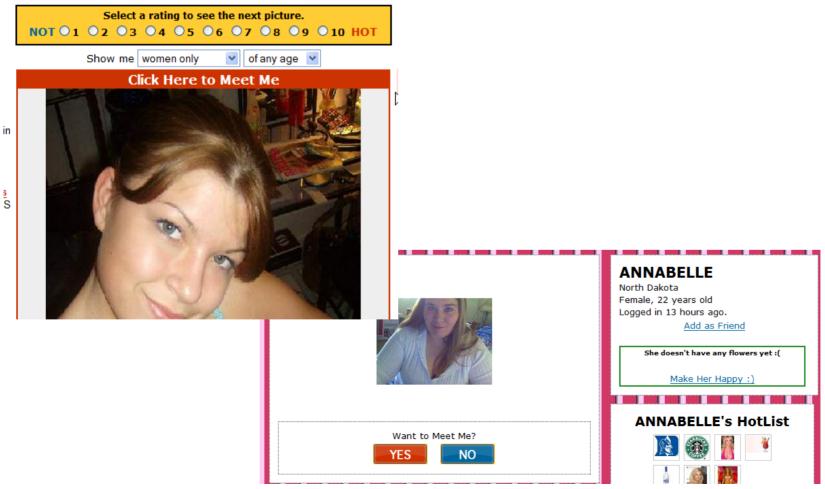


### • • Zillow Results

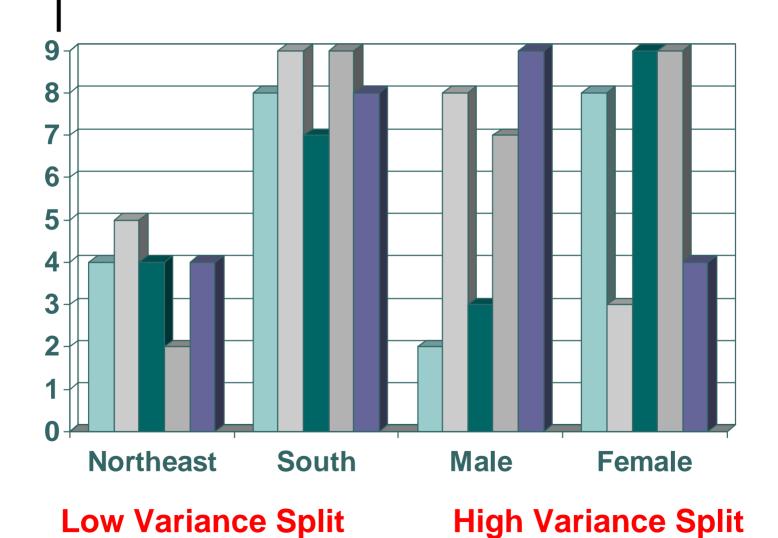


### Just for Fun... Hot or Not

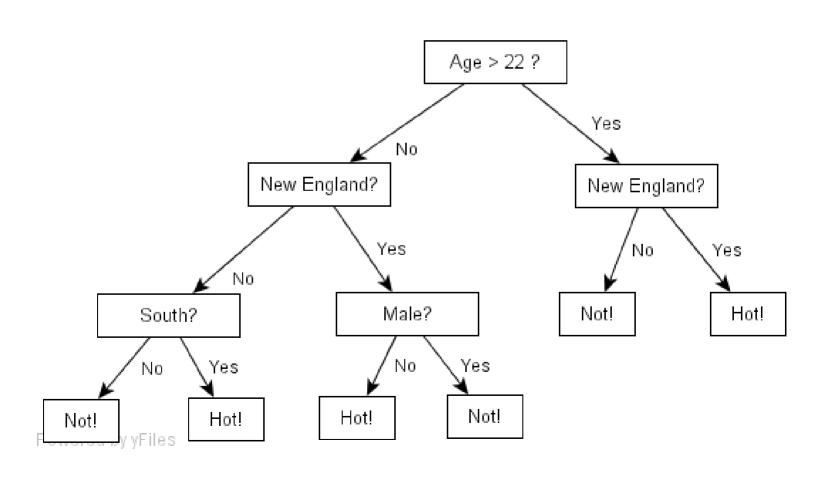
#### **HOT** or **NOT**.



#### Variance dividers



#### Just for Fun... Hot or Not



# Supervised and Unsupervised

- Clustering methods are unsupervised
  - There are no answers
  - Methods just characterize the data
  - Show interesting patterns
- Regression Trees are supervised
  - "answers" are in the dataset
  - Tree models predict answers

# • • Personal Ads

Anyone want to go to the Warriors exhib game tonight? - 33 (alamo squ

<u>s.f. bayarea craigslist &gt; women se</u>	eking men
all s.f. bayarea san francisco south	n bay east bay peninsula north bay
search for:	in: women seeking men
Poster's Age: min  [ Sun, 14 Oct 13:36:35 ] [fil	I don't ask for much - 23
Sun Oct 14  what are you doing? - 28 (laurel	Reply to: pers-448816051@craigslist.org Date: 2007-10-14, 10:59AM PDT
I wanna be a naughty girl! - 19	
restore my faith - 35 (concord / p Miss Being A Couple - 61	Hello there, to give you a little background I'm a fairly normal, quirky, fun, baseball-appreciating lady. I'm is bingo and reading books about social issues. No, I'm not 80, just a 23-year-old old soul type. I'm not high However, that does not mean I'm ok with being disrespected and unappreciated.
just sex with no strings - 26 (sar  When a girl misses you	I just want to meet a normal guy, and by normal I mean someone who doesn't drop acid to the point where me the first time he meets me or auto-assumes we're dating after ONE meeting. Seriously, guys?
Simply seeking 37 (san rafae paulina formerly of alameda - 36	- · · · · · · · · · · · · · · · · · · ·
SWF seeking SBM - 54 (San Fra	

## • • The Analysis

**Five Cities** 











**W4M Personal Ads** 

## • • Bayesian filter

If you listen to **NPR**, watch **Hardball**, and love the **Red Sox**, you may be the guy for me.

Please email me back.

I'm a professional with a grad school degree who has a sense of humor and loves the Sox.

	Sox	0.4
Boston	Red	0.35
<b>——</b>	Grad	0.2
	Professional	0.1
	Humor	0.1

# • • Bayesian filter

$$P(C | W) = P(C \& W) / P(W)$$

How often do the word and the city appear together?

How often does the word appear overall...

Rank these, and you have a list of the words most particular to a given city

### • Results

New York

Mets Pi

Lounges

Offense

**Desires** 

Musical

**Submissive** 

Create

Song

Oral

**Boston** 

Pink

Sox

Poetry

Intellectually

Punk

**Appreciation** 

**Exercise** 

Winter

**Education** 

Chicago

Cubs

Burbs

**Bears** 

Girlie

Insecure

Cheat

**Importance** 

**Blunt** 

Mouth

## • • Results

Los Angeles

Tee

**Excellent** 

Employment

San Francisco

Vegas

**Picnic** 

Meaningful

STD

Star

**Tasting** 

Lame

Hikes

Industry

French

Heat

.com

**Fitness** 

Kayaking

Entertainment

Cycling

Latino

#### Newsgroup Discussion



#### Home

New since last time: 5 messages

Description: Articles from dietandbody.com -- about diets and dieting, tips, tricks, FAQ, &A, recipes, exercises, diet reviews, weight loss success stories, and did-you-know. Join us and get lots of advice and support!

Almost-daily updates.



Discussions 9 of 184 messages view all »

Honest fat-burning tips- Workout for back muscles- Quinoa Salad With Garden Tomatoes- My own personal chef -- a success story

By diet\_and\_body - Oct 12 - 1 author - 0 replies

"Is cardio bad for me?" Correct link

By diet and body - Oct 11 - 1 author - 0 replies

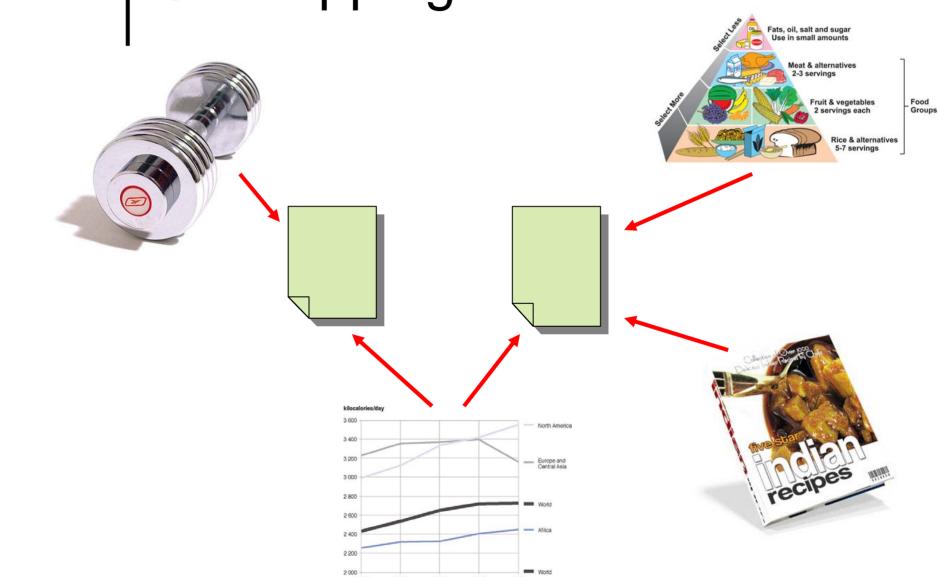
How to eat 5000 calories a day- Taste of Summer in Fall- Brown bag tips- Is cardio bad for me?

By diet and body - Oct 11 - 1 author - 0 replies

Chinese-Style Broccoli- Secrets of Gl Diet- Seven Foods that Cure- How to Stop Looking Fat

By diet\_and\_body - Oct 10 - 1 author - 0 replies

• • Overlapping themes



# Themes in a document

### alt.support.diet

#### How accurate are the calorie meters on gym equipment?

🖒 3 messages - Collapse all

#### McAlisters View profile

I use two different (readmills) one at home, the other at the YMCA. Both have the calories burned counters, the treadmill at the gym has imputs for age and weight, I'd assume for the calorie used math. I'm 58 years of age and weigh 205 and the machine tells me I burn about about 800 calories burned during an hours walk of 4.7 mph with an incline of 4.0 degrees. Does that sound about right for the calories burned?

Reply to author Forward Rate this post:

# Another word matrix

	Msg1	Msg2	Msg3	Msg4	Msg5
Gym	2	0	0	3	0
Calorie	0	2	1	1	3
Weigh	1	0	2	0	0
Carbs	0	3	0	0	2
Treadmill	1	0	0	2	0

**Actual Matrix** 

# • • Weights and features

	F1	F2	F3						
Gym	0	1	2		Msg1	M2	М3	M4	M5_
Calorie	2	0	1	F1	1	0	2	3	0
Weigh	2	2	1	<b>F</b> 2	0	2	1	1	3
		2	0	F3	1	0	2	0	0 ]
Carbs	1	Ü	3			\\\oig	ht Ma	triv	
Treadmill	0	1	2			vveig	iii ivia	UIA	

**Features Matrix** 

# • • Matrix factorization

	F1	F2	F3							
Gym	0	1	2			Msg1	M2	МЗ	M4	M5
Calorie	2	0	1		F1	<b>1</b>	0	2	3	0
Weigh	2	2	1	X	F2	0	2	1	1	3
Carbs	1	0	3		F3	1	0	2	0	0
Treadmill	0	1	2				Weig	ght M	atrix	

**Features Matrix** 

	Msg1	Msg2	Msg3	Msg4	Msg5
Gym	1	3	3	0	1
Calorie	0	2	4	1	3
Weigh	2	3	1	0	1
Carbs	0	1	1	0	2
Treadmill	3	2	0	2	2
	\				

**Current Guess** 

## • • Matrix factorization

Weight Matrix

	F1	F2	F3							
Gym	0	1	2			Msg1				
Gym Calorie	2	0	1		F1	<b>1</b>	0	2	3	0
Weigh Carbs	2	2	1	X	F2	0	2	1	1	3
Carbs	1	0	3		F3	1	0	2	0	0
Treadmill	0	1	2							

**Features Matrix** 

	Msg1	Msg2	Msg3	Msg4	Msg5
Gym	1	3	3	0	1
Calorie	0	2	4	1	3
Weigh	2	3	1	0	1
Carbs	0	1	1	0	2
Treadmill	3	2	0	2	2

**Current Guess** 

	Msg1	Msg2	Msg3	Msg4	Msg5
Gym	2	0	0	3	0
Calorie	0	2	1	1	3
Weigh	1	0	2	0	0
Carbs	0	3	0	0	2
Treadmill	1	0	0	2	0

Target Result

### • • Matrix factorization

Weight Matrix

	F1	F2	F3							
Gym	( 1	0	0			Msg1	M2	МЗ	M4	M5
Gym Calorie	0	1	1		F1	2	0	0	1	0
Weigh Carbs	0	0	2	X	F2	0	2	0	1	3
Carbs	0	1	0		F3	1	0	1	0	0
Treadmill	1	0	0							

**Features Matrix** 

	Msg1	Msg2	Msg3	Msg4	Msg5
Gym	2	0	0	3	0
Calorie	0	2	1	1	3
Weigh	1	0	2	0	0
Carbs	0	3	0	0	2
Treadmill	1	0	0	2	0

**Current Guess** 

	Msg1	Msg2	Msg3	Msg4	Msg5
Gym	2	0	0	3	0
Calorie	0	2	1	1	3
Weigh	1	0	2	0	0
Carbs	0	3	0	0	2
Treadmill	1	0	0	2	0

**Target Result** 

### • • Interpreting Features

	F1	F2	F3	
Gym	( 1	0	0	
Calorie	0	1	1	
Weigh	0	0	2	
Carbs	0	1	0	
Treadmill	1	0	0	

Theme 1	Theme 2	Theme 3
Gym	Calorie	Weigh
Treadmill	Carbs	Calorie

#### **Features Matrix**

	Msg1	M2	М3	M4	M5	
F1	2	0	0	1	0	
F2	0	2	0	1	3	
F3	1	0	1	0	0	<i>V</i>

Msg1 Msg2 Msg3 etc.
Theme 1 Theme 2 Theme 3
Theme 3

Weight Matrix

### "Diet and body" themes

**Atkins** 

Induction

South

Beach

Carbs

Chocolate

Black

Coffee

Olive

Broccoli

Gym

Weights

**Exercise** 

Running

Injured

Cook Recipe Fried

Home

Money

**Organic** 

Want

**Best** 

**Calories** 

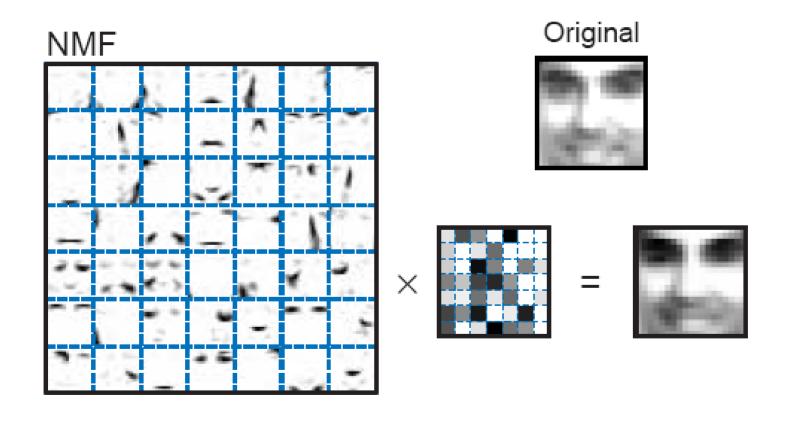
Weight

Fats

Protein

Cholesterol

### • • Side note: NMF for faces



## • • Methods covered

- Regression trees
- Hierarchical clustering
- k-means clustering
- Multidimensional scaling
- Bayesian Classifier
- Non-negative Matrix Factorization

# • • Other ideas

#### Finance

- Analysts already drowning in info
- Stories sometimes broken on blogs
- Message boards show sentiment

Extremely low signal-to-noise ratio

# • • Other ideas

- Product problems/ideas
  - Use support message boards
  - Extract themes
  - Understand recurring issues
  - Learn what features people want

# • • Other ideas

- Entertainment
  - How much buzz is a movie generating?
  - What psychographic profiles like this type of movie?

Of interest to studios and media investors