

Web Mining for Implicit User Affinities in On-line Communities

Matthew Smith
smitty@byu.edu
Brigham Young University
March 2005

Overview

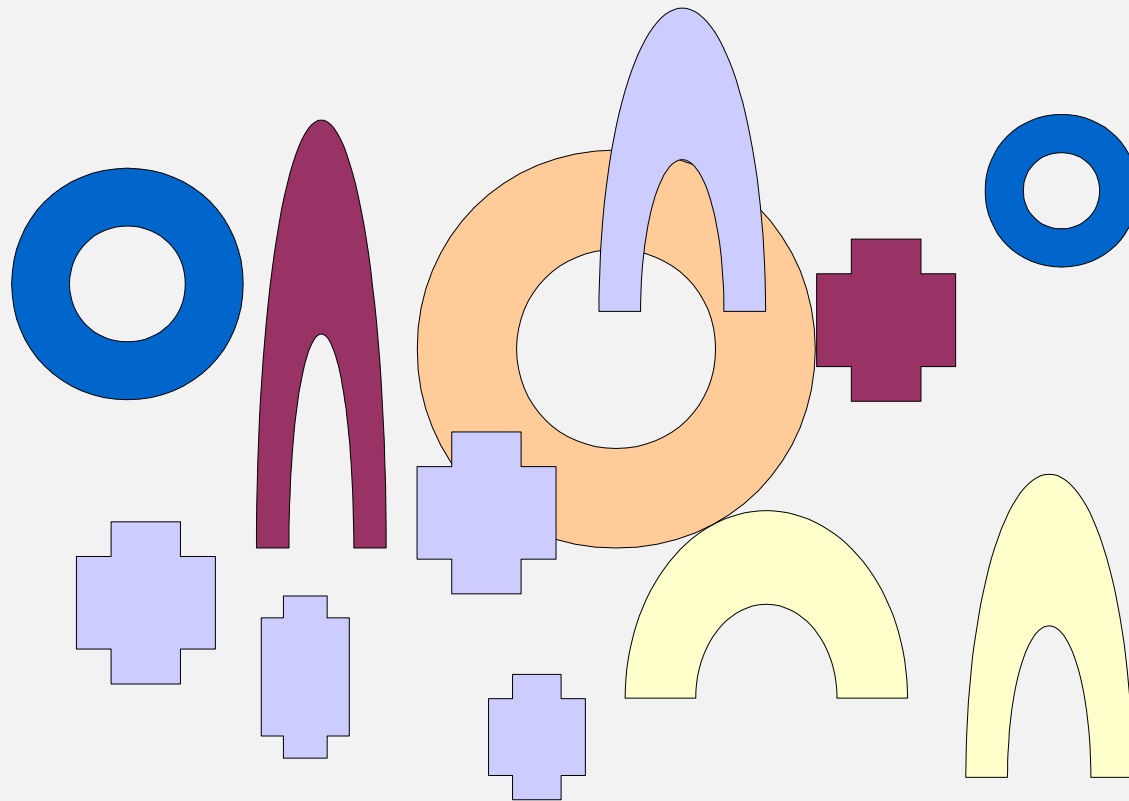
- Introduction
 - Affinity, Hypothesis, Scenario
- User Data
- Identifying Implicit User Affinities
- Quantifying Implicit User Affinities
- Conclusion

Affinity

- Definition:
 - An inherent similarity between persons or things
- Synonyms:
 - Relationship, connection, closeness, association, etc.

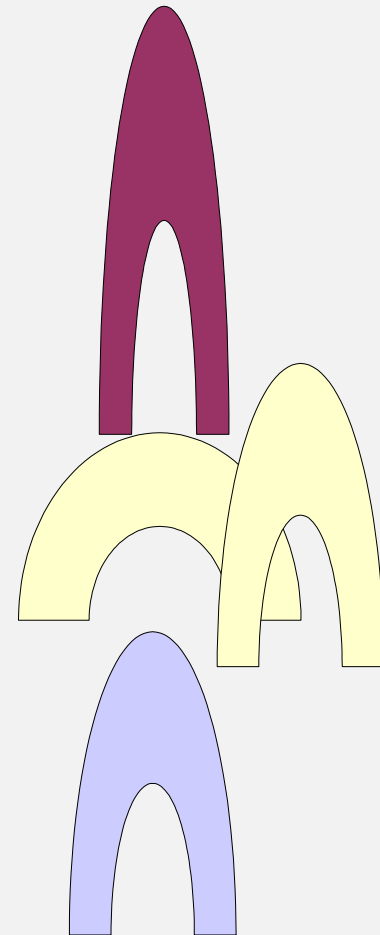
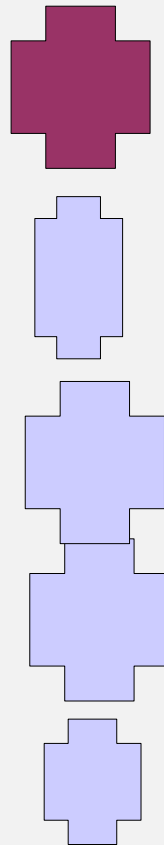
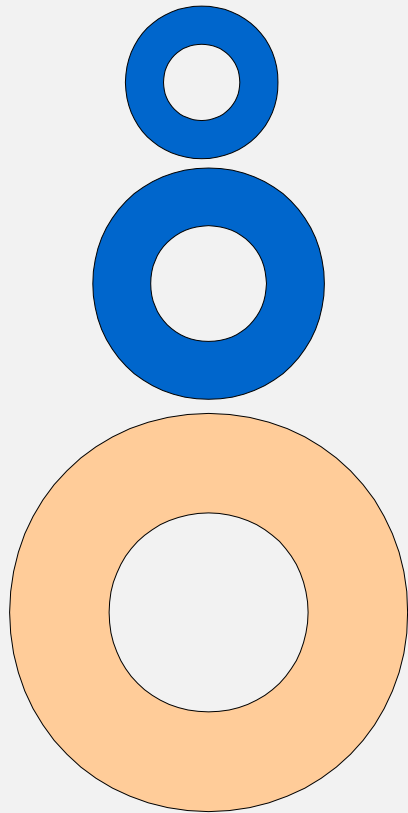
Basic Shapes

- What are some affinities?



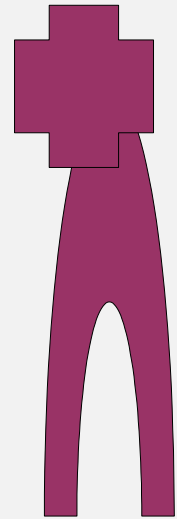
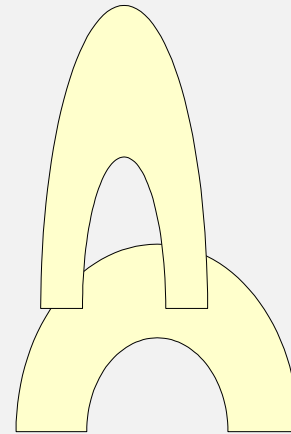
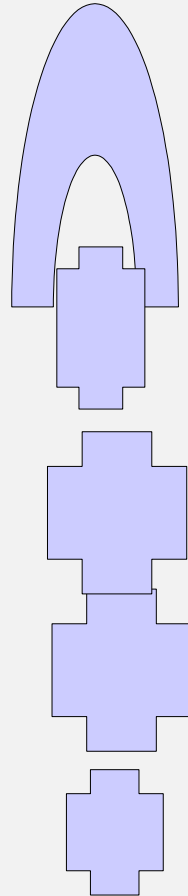
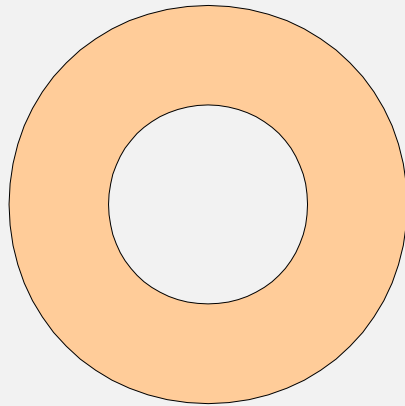
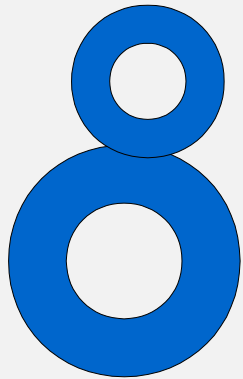
Clustered by Shape

- Three clusters [3 5 4]



Clustered by Color

- Five clusters [2 1 5 2 2]



Hypothesis

- Hypothesis:
 - Simple affinities can be discovered using basic user data within on-line communities
- Why might this be useful?
 - Better integrate new users into on-line communities
 - Interesting to learn how people within a community are related
 - Predict who users might like to meet

Use Case: Scenario

- You join an on-line community
 - *Pick one of your choice:* Google groups, Yahoo! Groups, **LDSMissions.com**, ACM, LDSMingle.com, Orkut, UUG, DevHood, uphpu, kddnuggets, or any other community you can think of.
- You want to know if you have any connections with anyone?
 - Who are those that you might know through other channels?
 - Who in the community do you share commonality (i.e., similar hobbies, similar geography, similar something, etc.) and who might you be interested in meeting?
 - Who of those similar to you is “well-connected” with other individuals you might like to meet?

User Data Collection

- Initial study: www.LDSMissions.com
 - Focus was on the basic data already being collected:
username, first name, last name, email address, city, state, zip, country, home ward & stake, and mission information (such as city, state/location, zip, country, areas served, companions, mission presidents, etc.)
 - Registration form <http://www.ldsmissions.com/us/?action=missionary.register>
 - What affinities might this data support?

LDS Missions Data

- Sample:

unique identifier

derived attribute

text mining required

from mission table

username	last name	first name	city	state	country	domain	mission presidents	areas served	companions	mission	mission country	mission location	start date	end date
			WJ	UT	UNITED STATES	msn.com	Pres Andrew Day	Taubate', Soracaba	Sis. Silva, Sis Nali	Brazil Sao Paulo North	Brazil	Brazil	Jul 1988	Jan 1990
			MUNTINLUPA CITY	PHILIPPINES	YAHOO.COM		PRESIDENT CAR	CATBALOGAN CITY	NGAWAKA, VALEN	Philippines Tacloban	Philippines	Philippines	May 1994	Mar 1996
			New Hartford	NY	UNITED STATES	email.byu.edu	Kradolfer	Alamogordo, Maran	Wardell, Staniszev	Arizona Tucson	United States	Arizona	Oct 2002	Apr 2004
			Elford	BC	CANADA	hotmail.com	Cook, Fillmore	North Philly, Potts	Bartlett, Nukaya, L	Pennsylvania Philadelphia	United States	Pennsylvania	Mar 2003	Mar 2005
			Orem	UT	UNITED STATES	hotmail.com	President Garrett,	Crescent View Ward	Elder Moore, Elde	Canada Calgary	Canada	Canada	Feb 2003	Mar 2005
			Bury, Lancashir		UNITED KINGDOM	yahoo.co.uk	Pres. Eldon. N. Tan	isle of white, salisbu	15.companions in	England London South	England	United Kingdom	Nov 1983	Sep 1985
			Pleasant Grove	UT	UNITED STATES	yahoo.com	Dewitt, Goodman,	Cosenza, Siracuz	Ellsworth, Greenh	Italy Catania	Italy	Italy	Nov 1994	Nov 1996
			Papatoetoe		NEW ZEALAND	hotmail.com	Pres. Wellis	Honolulu - Hawaii	Sis. Nishime, Sis. B	Hawaii Honolulu	United States	Hawaii	Sep 2002	Mar 2004
			Fairfax	VA	UNITED STATES	cox.net	Gary E. O'Brien	Biel, Ebnat Kappel,	Alan McLean, Mich	Switzerland Zurich	Switzerland	Switzerland	Nov 1974	Nov 1976
			Ogden	UT	UNITED STATES	yahoo.com	Christensen	Pdte. Roque Ensan	Williams, Enderle	Argentina Resistencia	Argentina	Argentina	Jan 1999	Jul 1999
			mesa	AZ	UNITED STATES	yahoo.com	ronald g. davis	madrid barrios:7,9,	flots	Spain Madrid	Spain	Spain	May 2003	May 2005
			Rexburg	ID	UNITED STATES	hotmail.com	President Hanks,	Irving sign, Irving Sp	Hermana Lyddon,	Texas Dallas	United States	Texas	Jun 1998	Dec 1999
			DeWinton	AB	CANADA	hotmail.com	Frank Bradshaw	Escindido, Carlsbad,	Miramar, Cardiff	California San Diego	United States	California	Jan 1973	Dec 1974
			Provo	UT	UNITED STATES	comcast.net	Hickman, Christia	Southport, Clinton, f	Golden, Loosli, Wi	North Carolina Raleigh	United States	North Carolina	Mar 1993	Feb 1995
			Omaha	NE	UNITED STATES	yahoo.com	H. Ray Hart	Metz, Boulogne sur	Mer, Cambrai, Bru	Belgium Brussels	Belgium	Belgium	Jun 1997	Jun 1999
			Bountiful	UT	UNITED STATES	mac.com	Sorensen, Thacke	Urbana, Springfield,	Johnson, Ohlson, J	Illinois Peoria	United States	Illinois	Jun 1997	Nov 1998
			Salt Lake City	UT	UNITED STATES	yahoo.com	Williams	La Roche sur Yon, f	Young, Hill, Grace	France Bordeaux	France	France	Jan 1999	Dec 2000
			salt lake city	UT	UNITED STATES	yahoo.com	Rod Tueller, Max	Sponchatoula, New C	Millar, Holt, Richar	Louisiana Baton Rouge	United States	Louisiana	Jan 1999	Dec 2000
			salt lake city	UT	UNITED STATES	yahoo.com	Parkers, Morenos	Las Delicias, Tarija,	Matsumura (steph	Bolivia Cochabamba	Bolivia	Bolivia	Jan 1999	Jul 2000
			antioch	CA	UNITED STATES	aol.com				Venezuela Valencia	Venezuela	Venezuela	Apr 1999	Mar 2001
			Kaysville	UT	UNITED STATES	msn.com	Schreiber		Watts, Schmidt, Pl	Germany Hamburg	Germany	Germany	Jan 1981	Aug 1982
			The Woodlands	TX	UNITED STATES	vogtengineering	J. Ballard Washbu	Prescott, Phoenix, T	Suzette Whiting, H	Arizona Phoenix	United States	Arizona	Jan 1988	Aug 1989
			Taylorsville	UT	UNITED STATES	hotmail.com	President Bernard	Somerset, George	Pann, Christense	Kentucky Louisville	United States	Kentucky	Dec 2002	Dec 2004
			Norman	OK	UNITED STATES	hotmail.com	Marshal, Cahoon	Veracruz, Puebla, A	Northcutt, Borget,	Mexico Veracruz	Mexico	Mexico	Apr 1979	Apr 1981
			CANYON COUN	CA	UNITED STATES	YAHOO.COM	PRESIDENT LEE	HONDURAS TEGUCIGALPA		Honduras Tegucigalpa	Honduras	Honduras	Jan 2004	Aug 2005
			Cedar Rapids	IA	UNITED STATES	mchsi.com	Gosta Berling, Joh	Bergen, Sarpsborg,	Gary Stoddard, Ste	Norway Oslo	Norway	Norway	Jan 1975	Jan 1977
			Sandy	OR	UNITED STATES	netzero.com	Hobbs, Hunter			New Hampshire Manches	United States	New Hamps	Feb 1997	Feb 1999
			los angeles	CA	PHILIPPINES	yahoo.com.	Daniel rogers	kansas city	sister burgess	Missouri Independence	United States	Missouri	Oct 1998	Sep 2000
			Boise	ID	UNITED STATES	yahoo.com	Donald Hinton, Te	Aberdeen, Ngau Tai	Ure, J. Chan, Vale	China Hong Kong	China	China	Sep 2002	Oct 2004
			Brawley	CA	UNITED STATES	juno.com	Paul Felt	White Cone, Twin L	Phil Smith, E. Rog	Arizona Mesa	United States	Arizona	Jan 1970	Jan 1972

Hidden to preserve anonymity

Identifying Implicit Affinities

- How can affinities be implied?
- Method: **Similarity Clustering and Aggregation**
 - Names: first, last, middle, maiden
 - Mission: name, presidents, companions
 - Geographically
 - Mission: region, country, state/location, city, zip, areas served
 - Home: region, country, state, city, zip, home stake, home ward
 - Email address domains (i.e., hotmail.com, yahoo.com, msn.com, byu.edu, etc.)

Data Aggregation

- Email Domains

<u>domain</u>	<u>count</u>		
hotmail.com	4863	myldsmail.net	59
yahoo.com	2316	worldnet.att.net	59
aol.com	1413	attbi.com	53
juno.com	687	comcast.net	50
email.byu.edu	500	latinmail.com	49
msn.com	302	byuh.edu	49
earthlink.net	124	uswest.net	46
cc.usu.edu	115	prodigy.net	44
byu.edu	101	mstar2.net	41
cs.com	97	sisna.com	39
excite.com	93	infowest.com	38
usa.net	77	weber.edu	35
netzero.net	68		
netscape.net	67		
home.com	65		
cox.net	62		
byui.edu	62		

Data Aggregation

- First names

<u>first_name</u>	<u>count</u>
David	385
Michael	284
Ryan	224
Scott	209
Jason	200
John	181
Daniel	174
Jared	167
Matthew	160
Robert	149
Brian	147
Mark	146
James	139
Aaron	132
Nathan	125
Eric	125
Benjamin	124
Adam	124
Kevin	121
Chris	116
Brandon	115
Richard	112
Paul	111
Andrew	110
Joshua	110
Steven	109
Matt	104
Jeremy	104
Christopher	101
Justin	99
...	
Burdette	1

- Last names

<u>last_name</u>	<u>count</u>
Smith	168
Johnson	126
Jones	104
Anderson	89
Brown	82
Christensen	74
Jensen	74
Hansen	72
Nelson	67
Peterson	66
Taylor	61
Clark	56
Williams	55
Davis	53
Allen	49
Thompson	48
Miller	48
Larsen	47
Wilson	45
Wright	45
Adams	42
Hall	41
Harris	40
Hill	39
White	38
Walker	36
Nielsen	33
Hatch	33
Sorensen	32
Olsen	32
...	
Pixton	1

Data Aggregation

- Geography (Country, State, City)

<u>country</u>	<u>count</u>
United States	14272
Canada	320
Philippines	256
Mexico	213
Brazil	155
Chile	111
Australia	110
United Kingdom	110
Argentina	95
New Zealand	79
Peru	69
Spain	33
Guatemala	26
Colombia	23
France	21
Ecuador	21
Bolivia	20

Germany	20
Uruguay	18
Venezuela	18
Sweden	16
Netherlands	15
Puerto Rico (US)	15
Portugal	13
Switzerland	13
Italy	11
Costa Rica	11

<u>state</u>	<u>count</u>
UT	6223
CA	1654
ID	1154
AZ	1116
WA	552
TX	478
NV	316
OR	305
CO	245
HI	209
ALB	193
FL	189
VA	135
GA	122
NM	109
WY	105
MO	92

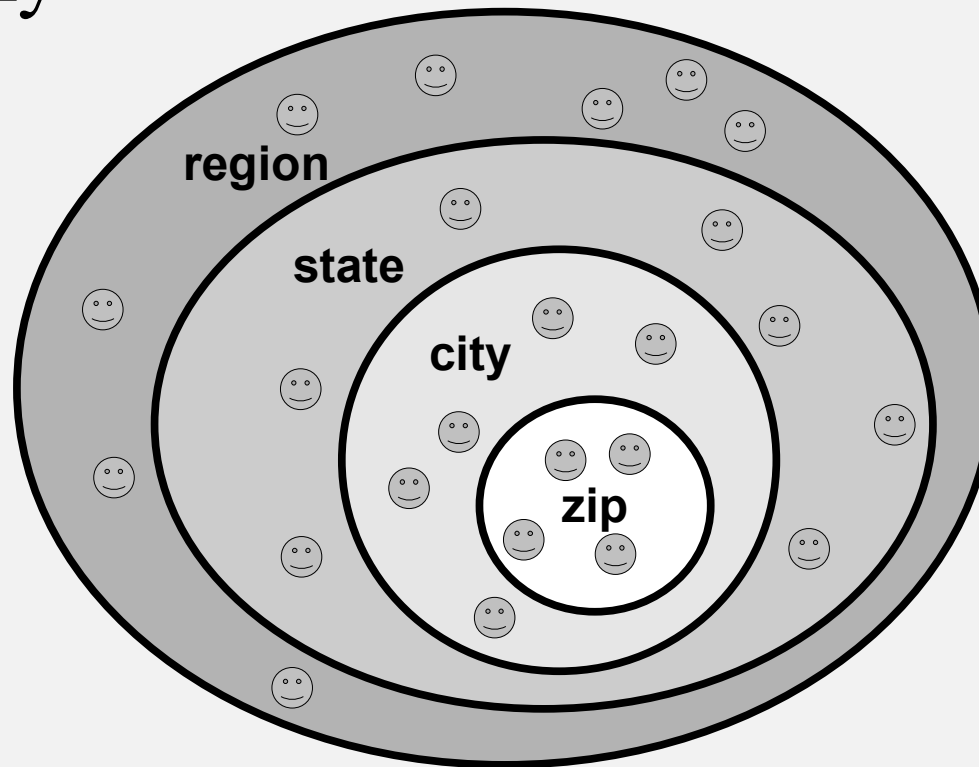
NC	90
AL	87
OH	87
IL	85
TN	81
IN	76
NY	76
MI	65
BC	64
MD	61
AK	58
KS	57

<u>city</u>	<u>count</u>
Provo	751
Salt Lake City	566
Orem	479
Mesa	402
Sandy	297
Ogden	189
LOGAN	175
Layton	167
Bountiful	165
Las Vegas	164
West Jordan	157
Idaho Falls	146
Rexburg	132
St. George	121
pocatello	111
Taylorsville	109
Boise	107
Gilbert	103
West Valley City	102
SLC	98
Phoenix	98
Murray	98
laie	97
Cedar City	91
South Jordan	86
Kaysville	86
Farmington	84
Draper	77
Springville	72

Visualization: Venn Diagram

(stacked)

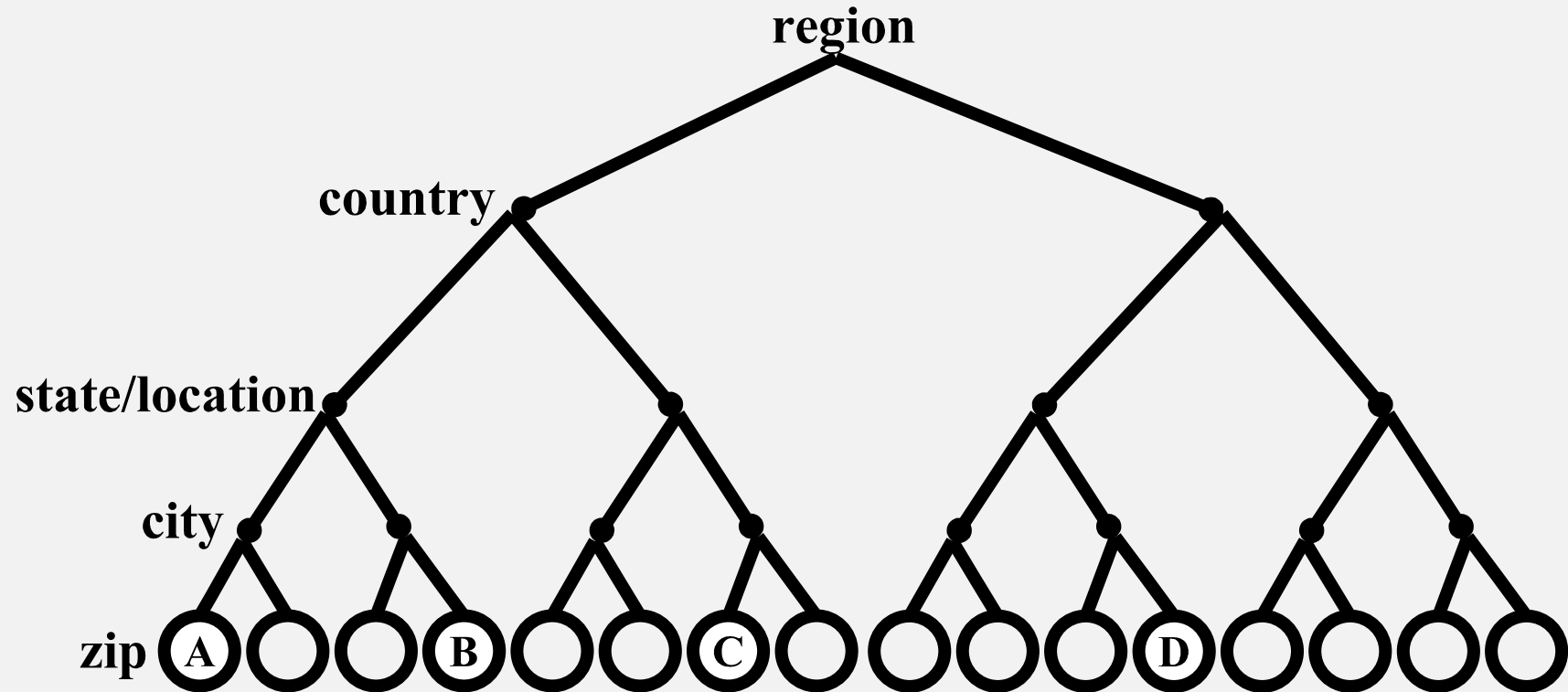
- Geography



Which group is most interesting?

Visualization: Attribute Hierarchy

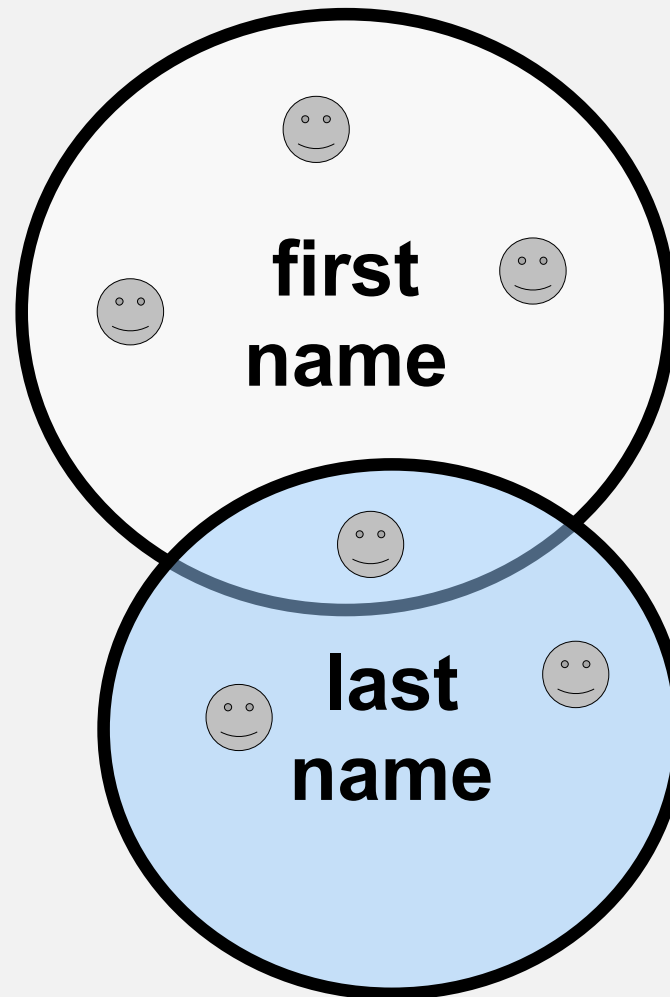
- Geographical (related attributes)



Visualization: Venn Diagrams

- Names

- first name match
- last name match
- both match



LDS Missions Example

- Login to website...
 - <http://www.ldsmissions.com>

Quantifying Implicit Affinities

- Why? Important for comparing affinities
- All affinities are not equal
- Affinities are valued differently among users
- Affinities can be aggregated to form a combined affinity score
 - In some cases, quantification will be most useful as a combined calculation

Approach

- Weight affinities equally
 - easiest to implement but less accurate
- Survey all users, then weight affinities accordingly
 - Result may be applied either:
 - Individually
 - Customized affinity scores for each user
 - Collaboratively
 - Average of all user surveys
- Weight affinities by an expert
- Other ideas?

Conclusion

- Useful affinities can be discovered implicitly
 - Using already collected data.
- It is a novel approach to auditing your data
- It is a small world after all.
 - That is, most people are connected through some sort of affinity.
 - All people are are connected through usually small affinity chains.