

Machine Learning with R using Qualitative Research Codes

Marla Stuart, PhD, MSW

Fellow, Berkeley Institute for Data Science

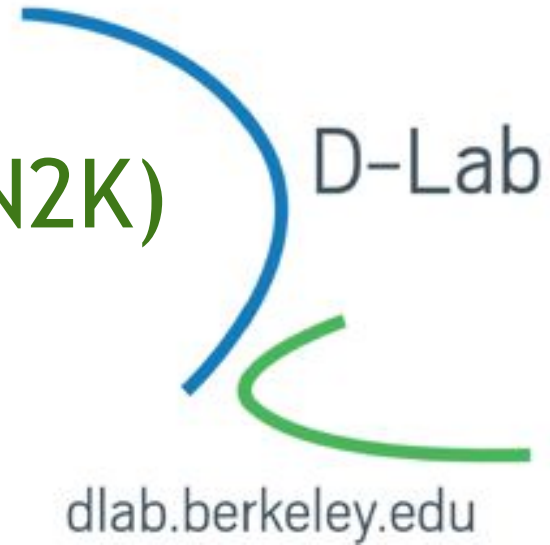
Fellow, Guizhou Berkeley Big Data Innovation Research Center

Josué Meléndez Rodríguez, MA, MSW

Qualitative Research Lead & Senior Data Science Fellow, D-Lab



- It's Okay Not to Know (IOKN2K)
- 280 workshops
- 1,100 consultations
- working groups
- special research projects
- approx. 6,000 scholars served per year



Agenda

- Introductions
- Qualitative Research with QDA Software (Josué)
 - Qualitative Coding & Analysis
 - Data in Sample Project
 - Qualitative Data Analysis Software
- Qualitative Research with Machine Learning (Marla)
 - The case for using machine learning
 - What is Machine Learning?
 - Unsupervised ML for theming
 - An Example
- Questions and Discussion

Qualitative Research with QDA Software

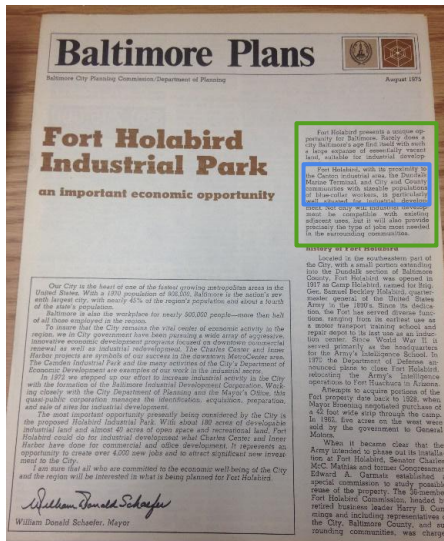


Qualitative Research...

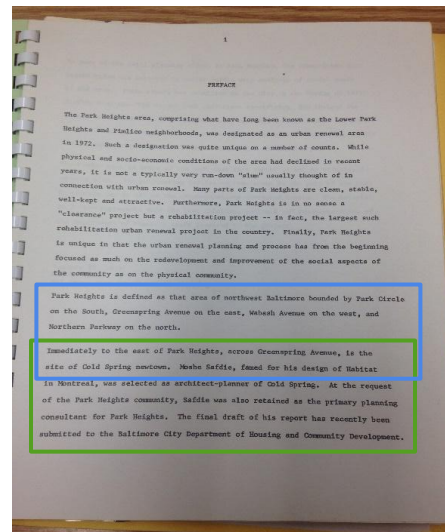
What are Codes? What is Coding?

Coding is a way of organizing the data around some common idea, concept, or category ACROSS sources.

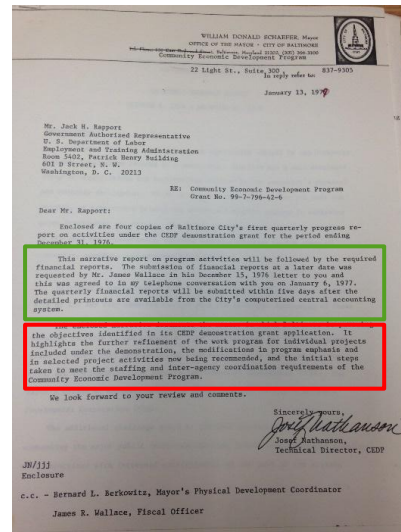
A



B



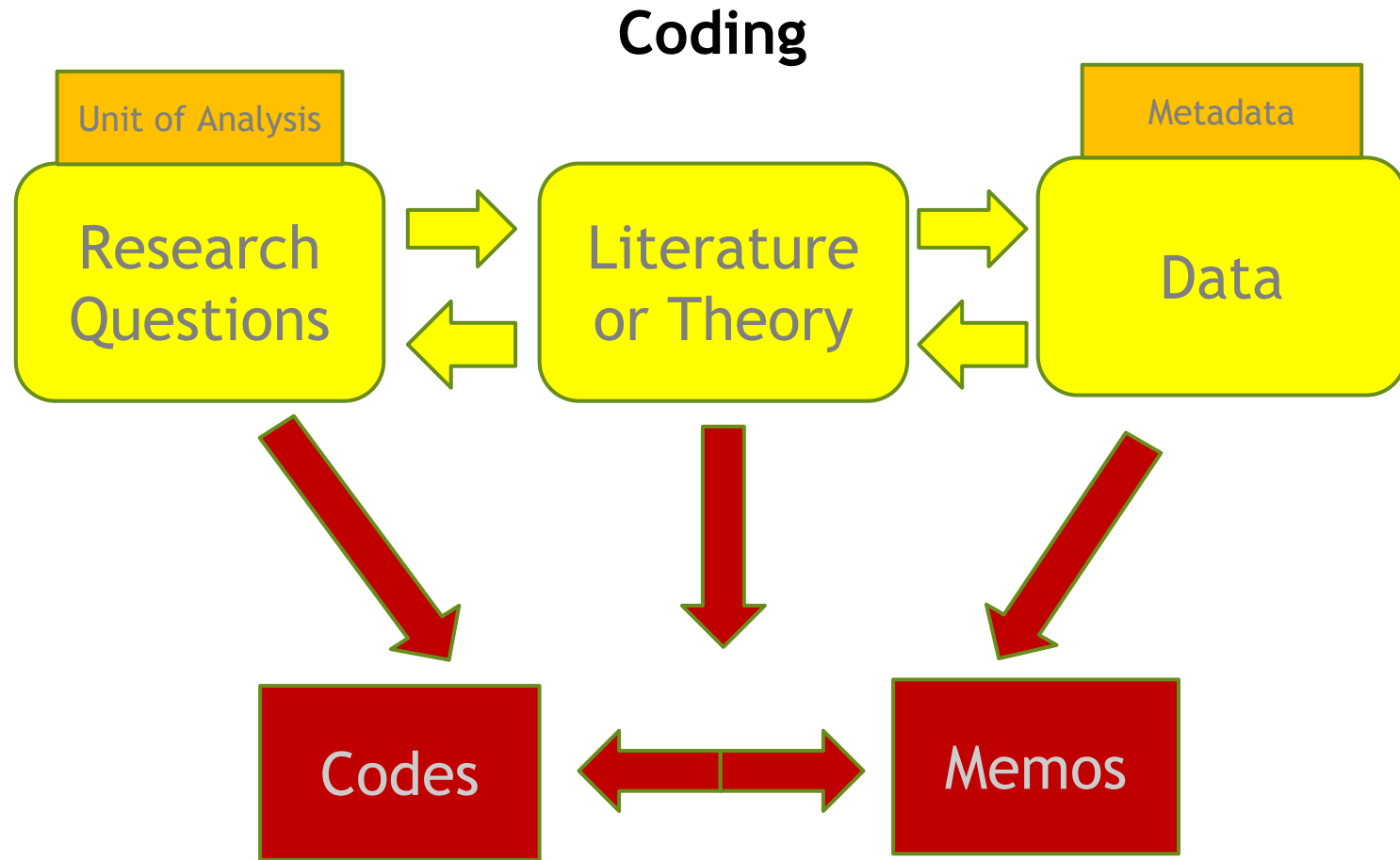
C



The code of "financial planning" is applied to the selected text from documents A, B, and C, because they all discuss this topic.

Qualitative Research...

Coding as a Multistep Nonlinear Process



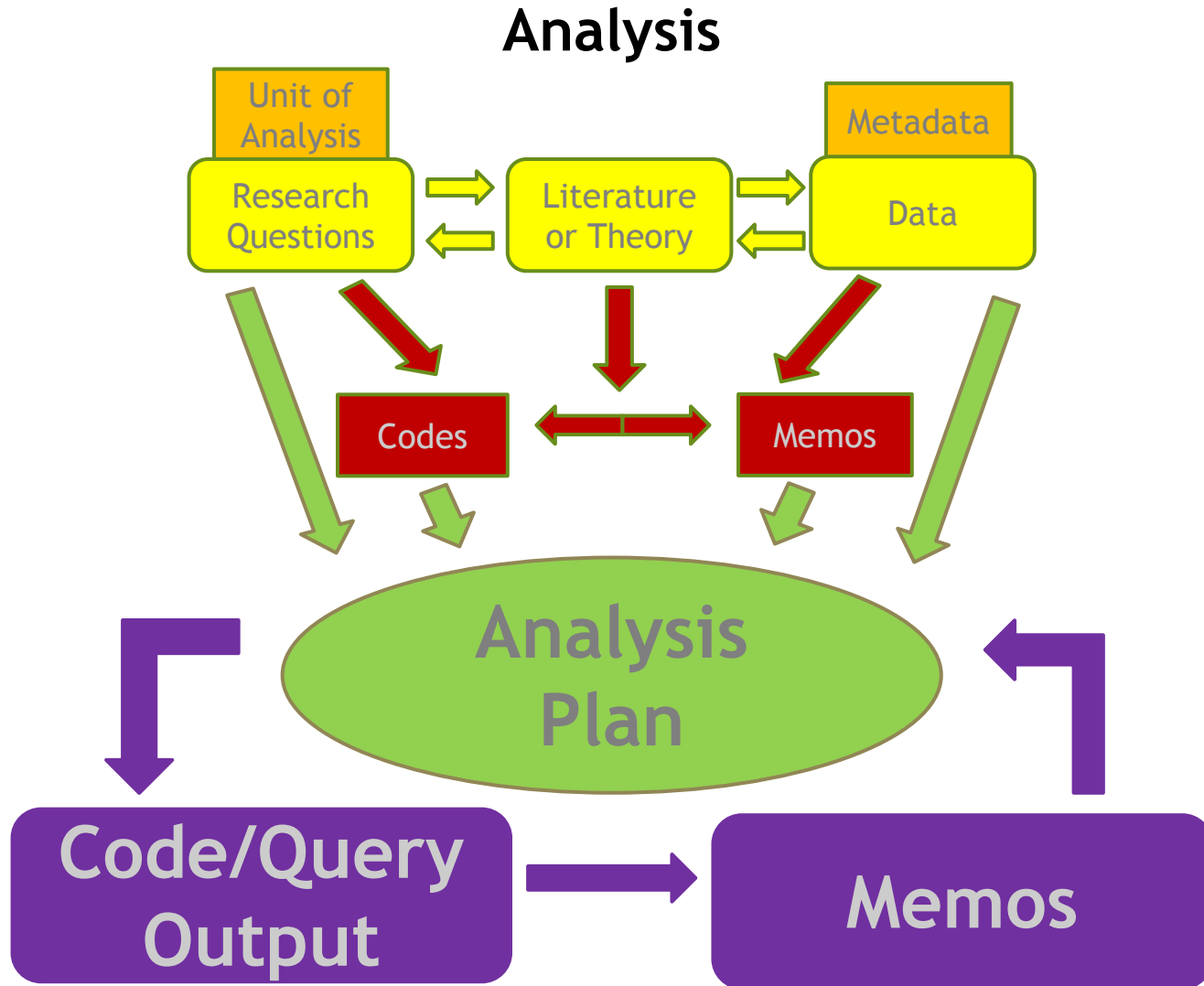
What is Analysis?

The process of identifying themes related to your research findings. This is different than identifying ideas/concepts/topics that come up throughout your data set. It's "bigger picture" stuff...

- Overarching Themes
 - What is happening in your data overall?
- Subgroup Themes
 - What is happening in your data for specific subgroups?
- Typology Themes
 - What is happening in your data by specific dimensions of coded data?

Qualitative Research...

Analysis as a Multistep Nonlinear Process



Qualitative Research...

Data Used in Sample Project

Data Sources

journal articles focusing on intersectionality and social work practice

- ▶ Intersectionality as a Useful Tool: Anti-Oppressive SW and Critical Reflection
- ▶ Children's Health Disparities as Embodiment of Social Class
- ▶ Recognizing Students of Color as Holders and Creators of Knowledge
- ▶ Undocumented Mexican Mothers in the Current Policy Context
- ▶ Identity, Oppression, and Power: Feminisms and Intersectionality Theory
- ▶ African American Women's Perceptions of Intersection of DV & HIV/AIDS

Codes

- ▶ Intersectionality
- ▶ Practice Level
 - ▶ Micro
 - ▶ Macro
- ▶ Privilege and Oppression
 - ▶ Sexism
 - ▶ Racism
 - ▶ Classism
 - ▶ Heteronormativity
 - ▶ Cisgenderism
 - ▶ Ableism
 - ▶ ...
- ▶ Practice Recommendations

QDA Software as a Tool for Coding & Analysis

What It Does

- Structure & Organize
- Code & Retrieve
- Memo
- Explore
- Query
- Visualize

What It Does Not

- Error-Free Auto Coding
- Analytic Thinking
- Eliminate Bias
- Advanced Quantitative Analysis
 - MAXQDA engages in some quantitative analysis

Potential Benefits

- Frees Time to Focus on Analysis
- Can Deal with Large Data Sets
- Improves Auditability
(among some audiences)
- Improves Credibility
(among some audiences)

Potential Drawbacks

- Requires Learning the Software
- May Force Production of Meaningless Findings
- May Create Pressure to Engage Excessive Features

Qualitative Research...

Overview of Relevant QDA Software Programs



QDA Miner



HyperResearch

ANSWR



Transana

Aquad (open source)



Quirkos (visual exports)

Saturate (app-based)



Qualitative Research and Machine Learning

The background of the slide is white with abstract green geometric shapes on the right side. These shapes are overlapping triangles and polygons in various shades of green, ranging from light lime to dark forest green. The shapes are positioned on the right side of the slide, creating a modern, tech-oriented aesthetic.

The case for machine assisted text analysis



The Qualitative Report

Volume 16 | Number 3

Article 6

5-1-2011

Compatibility between Text Mining and Qualitative Research in the Perspectives of Grounded Theory, Content Analysis, and Reliability

Chong Ho Yu
Arizona State University, chonghoyu@gmail.com

Angel Jannasch-Pennell
Arizona State University, angel@asu.edu

Samuel DiGangi
Arizona State University, sam@asu.edu

99
citations

Benefits of Machine Learning

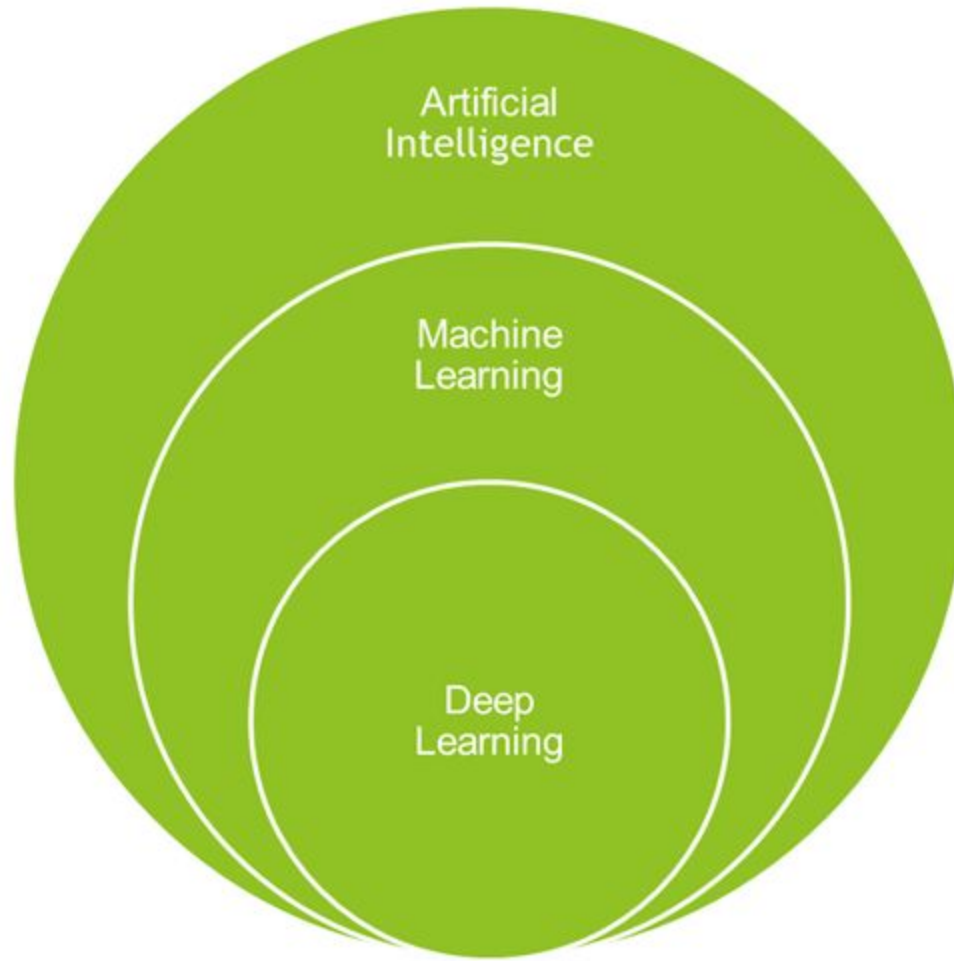
Natural extension of Content Analysis and Grounded Theory (data-informed)

Facilitates use of more text

Replicable
Transparent
Fit statistics
More reliable

Efficient

What is Machine Learning?



AI: Techniques that enables machines to mimic human logic, make predictions, and act on those predictions.

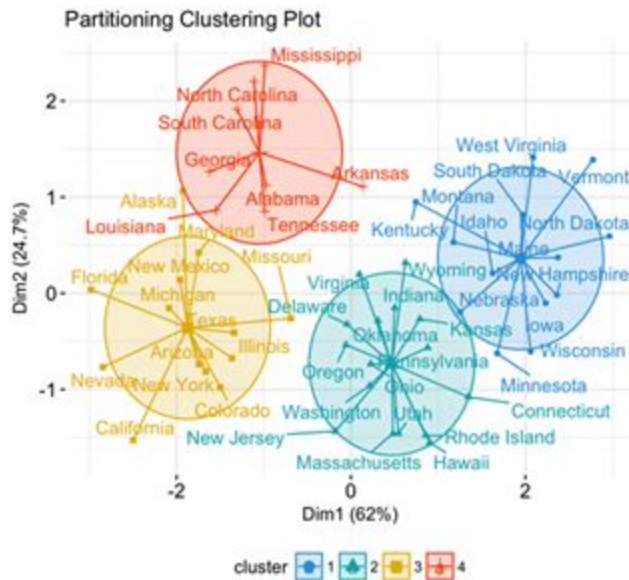
ML: Techniques by which a computer deploys user-defined logical sequences and statistical methods to predict an observed or unobserved class. If predictions are inaccurate, adaptations to the rules are user-defined.

DL: Techniques by which a computer deploys user-defined logical sequences and statistical methods to predict an observed or unobserved class. If predictions are inaccurate, adaptations to the rules are computer generated.

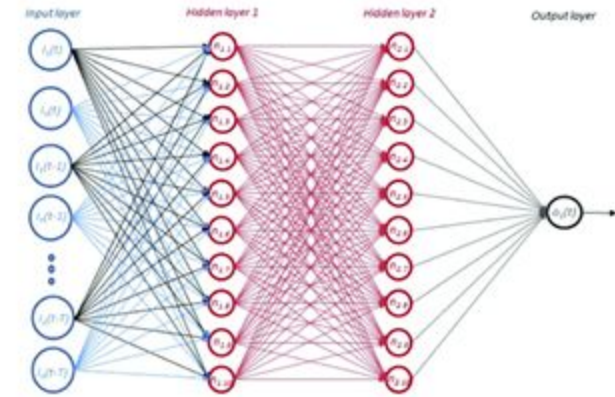
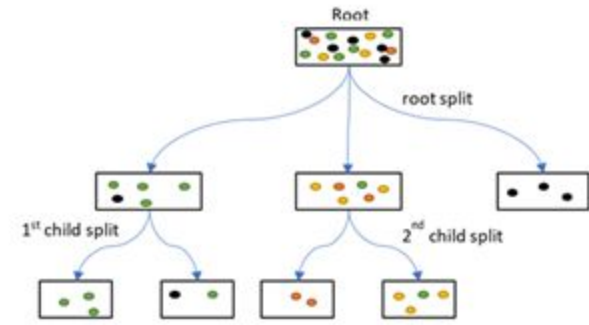
Equations and Algorithms

	Equations	Algorithms
	Column-based Relationships among variables relative to all observations Measures of uncertainty (model fit) well developed Statistics	Row-based Relationships among observations relative to all variables Measures of uncertainty less well developed Computer Science/Engineering
Predicting observed outcomes	Regressions	Supervised Machine Learning
Predicting unobserved outcomes	Factor Analysis Latent Class Analysis Structural Equation Modeling	Unsupervised Machine Learning

Common Machine Learning Methods



Unsupervised Models
Predicting an unobserved outcome



Supervised Models
Predicting an observed outcome



Unsupervised Machine Learning for qualitative code analysis / theming

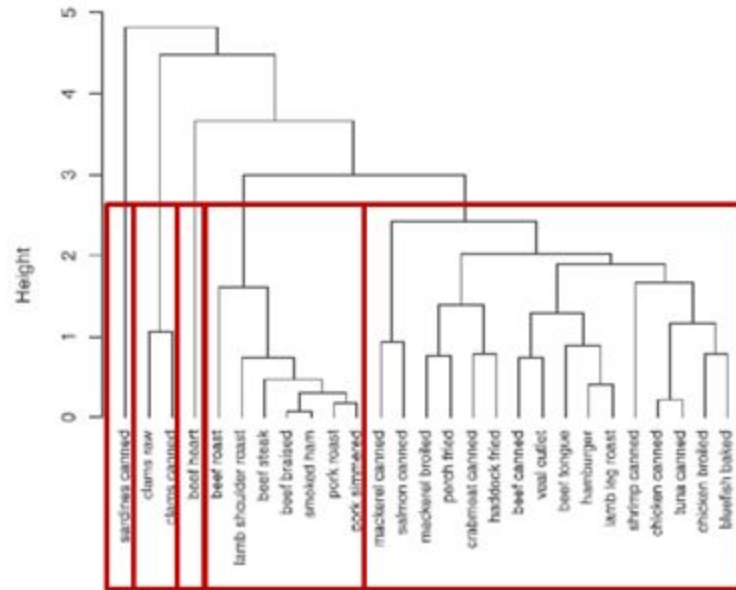
clustering: distance measure

$$d_{ji} = \sqrt{\sum_{p=1}^p (x_{ip} - x_{jp})^2}$$

observation	energy	protein	fat	calcium	iron
BEEF BRAISED	340	20	28	9	2.6
HAMBURGER	245	21	17	9	2.7
BEEF ROAST	420	15	39	7	2.0
BEEF STEAK	375	19	32	9	2.6

$$d = \sqrt{(340 - 245)^2 + (20 - 21)^2 + (28 - 17)^2 + (9 - 9)^2 + (2.6 - 2.7)^2} = 95.64$$

hierarchical



cluster	energy	protein	fat	calcium	iron
1	340	19	29	9	2.50
2	170	20	8	13	1.45
3	160	26	5	14	5.9
4	57	9	1	78	5.70
5	180	22	9	367	2.5

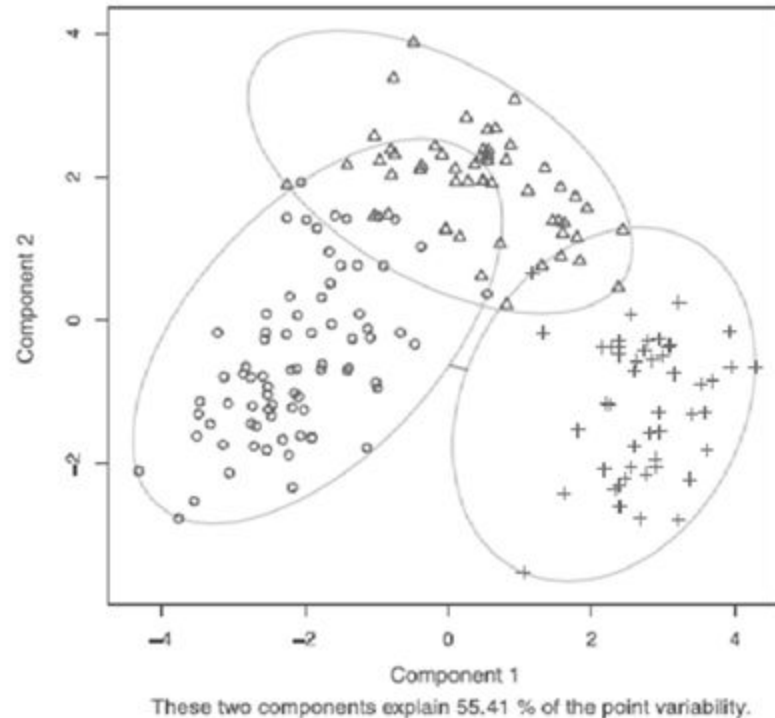
Logic

1. Each case or observation starts as its own cluster.
2. Clusters are combined two at a time until all clusters are merged into a single cluster.

Algorithm

1. Define each observation (row, case) as a cluster.
2. Calculate the distances between every cluster and every other cluster.
3. Combine the two clusters that have the smallest distance. This reduces the number of clusters by one.
4. Repeat steps 2 and 3 until all clusters have been merged into a single cluster containing all observations.

partitioning: k-means and PAM



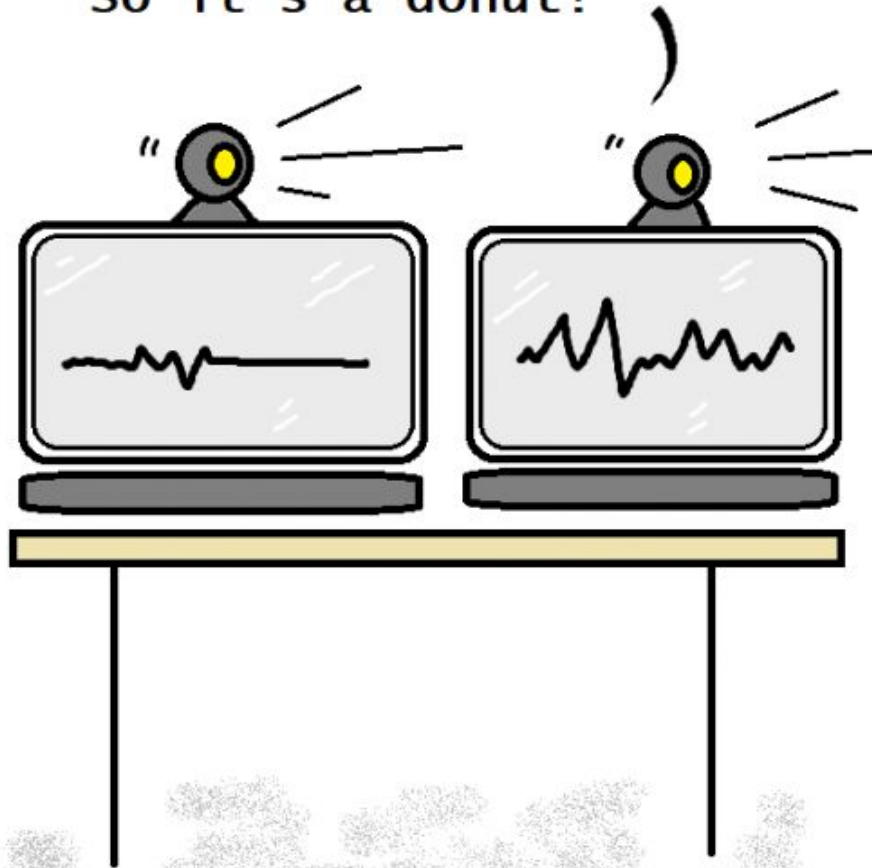
Logic

1. Observations divided into K groups.
2. Observations are reshuffled to form the most cohesive clusters possible.

Algorithm

1. Select K centroids (K rows chosen at random).
2. Assign each data point to its closest centroid.
3. Recalculate the centroids as the average of all data points in a cluster.
4. Assign data points to their closest centroids.
5. Continue steps 3 and 4 until observations aren't reassigned or the maximum number of iterations is reached.

"It's motionless, soft
and brightly colored.
So it's a donut!"



Akell...



[Linkedin.com/in/azi-azimi/](https://www.linkedin.com/in/azi-azimi/) © AziAzimi

An Application to Josué's Sample Study

Data (codes)

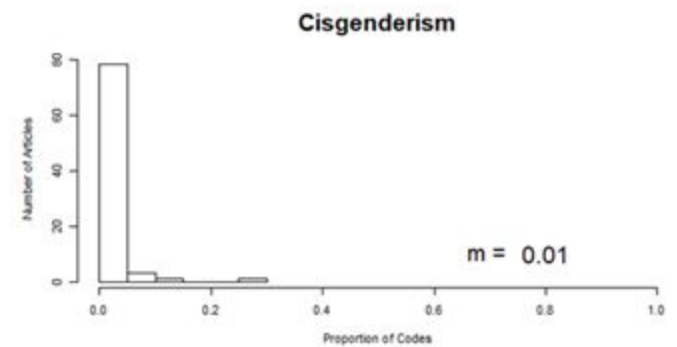
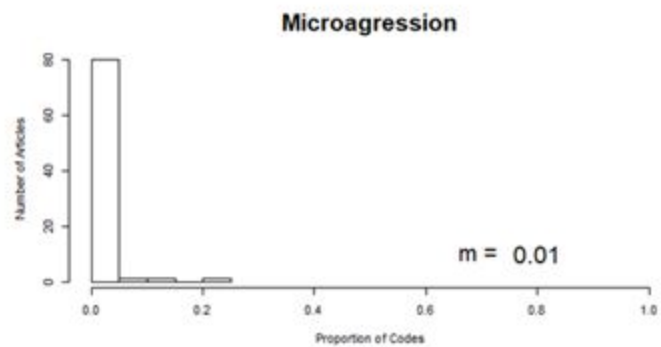
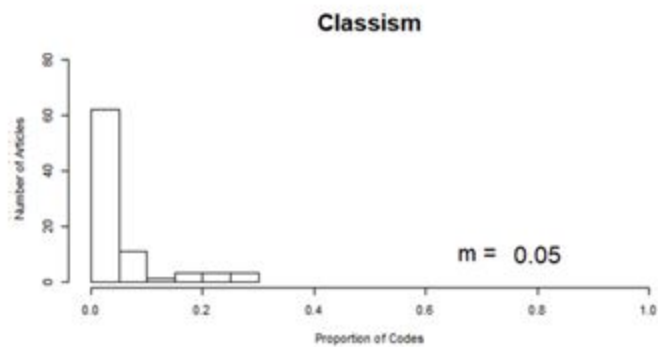
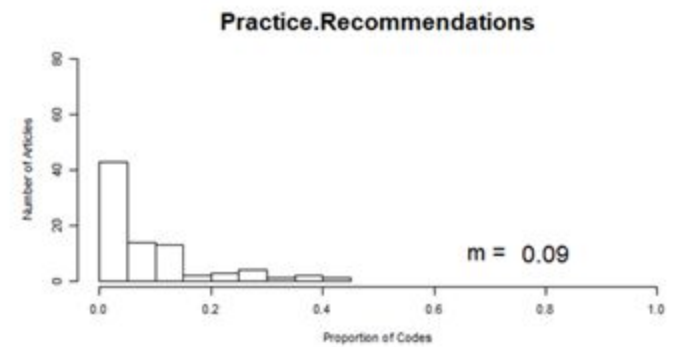
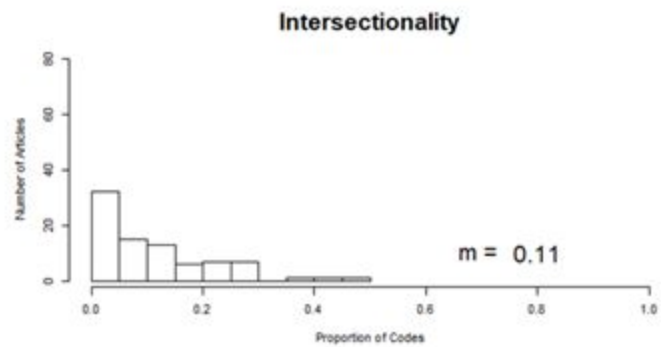
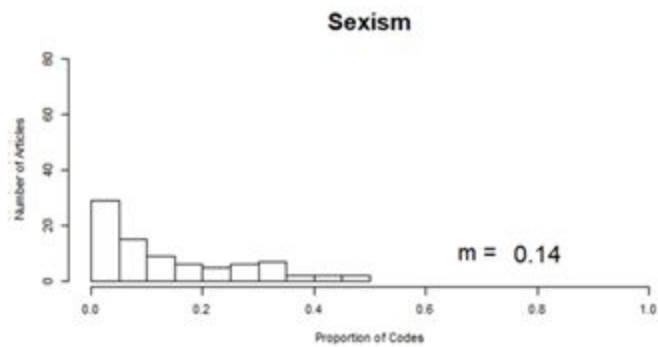
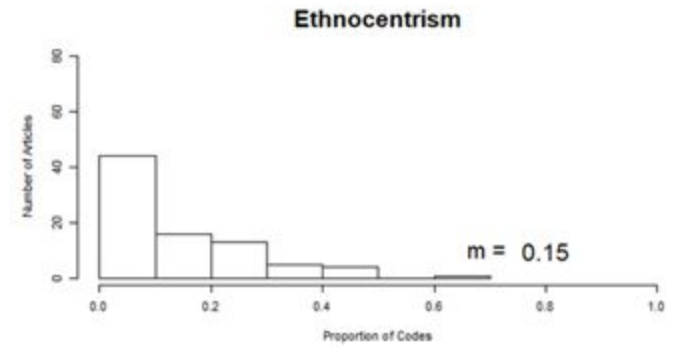
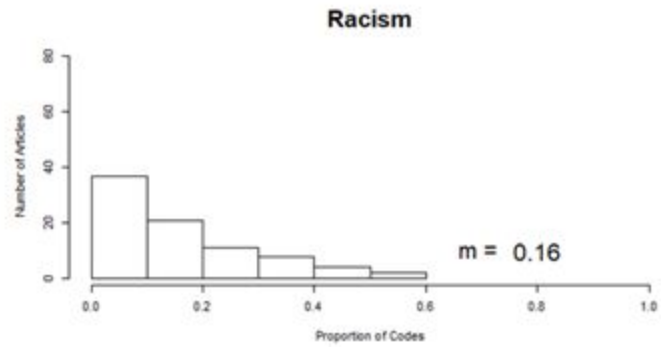
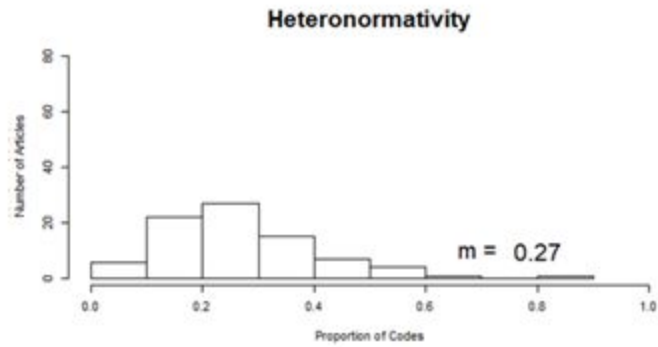
josue.csv - Excel

File Home Insert Page Layout Formulas Data Review View Developer Acrobat Tell me what you want to do...

E9 X ✓ fx 0

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	Document	Segments	Intersectionality	Macro	Micro	Microaggression	Ableism	Cisgenderism	Ethnocentrism	Heteronormativ	Classism	Sexism	Racism	Practice	
2	1	238	23	0	13	13	5	9	17	100	5	13	35	5	
3	2	116	5	2	4	4	0	4	9	31	1	13	38	5	
4	3	28	8	0	0	0	0	0	2	2	1	9	4	2	
5	4	67	1	3	1	0	0	0	27	7	1	2	25	0	
6	5	193	19	1	0	0	0	0	51	21	4	64	18	15	
7	6	69	0	0	0	0	0	0	7	15	0	1	36	10	
8	7	73	8	0	1	0	0	0	7	26	3	5	21	2	
9	8	95	5	0	0	0	3	0	3	25	15	33	10	1	
10	9	130	2	2	1	0	0	0	33	21	20	37	13	1	
11	10	57	3	0	0	0	0	0	13	18	1	5	14	3	
12	11	131	6	0	1	1	0	0	34	28	0	16	23	22	
13	12	36	0	1	1	0	0	0	8	8	0	3	13	2	
14	13	102	23	1	2	2	0	0	11	28	2	3	6	24	
15	14	97	12	0	3	3	1	0	9	23	2	2	7	35	
16	15	110	1	0	0	0	0	0	52	22	6	11	5	13	
17	16	91	5	1	0	0	0	0	13	23	1	39	5	4	
18	17	41	2	0	0	0	0	0	7	13	1	0	15	3	
19	18	127	2	0	5	0	0	0	35	23	3	2	55	2	
20	19	64	11	0	0	0	0	0	9	20	0	17	4	3	
21	20	167	3	0	1	0	0	0	4	88	2	52	10	7	
22	21	167	3	5	13	0	2	0	18	38	10	0	16	62	
23	22	106	2	0	0	0	0	0	5	36	2	20	6	35	
24	23	89	21	1	1	0	0	0	4	24	1	16	10	11	
25	24	159	4	1	1	0	0	0	4	36	1	42	67	3	
26	25	27	0	0	0	0	0	0	1	22	0	0	0	4	
27	26	205	8	3	0	0	2	0	24	104	9	44	6	5	
28	27	136	29	1	0	0	0	0	6	31	6	44	4	15	
29	28	155	9	0	0	0	0	0	14	42	4	72	1	13	
30	29	127	3	2	10	0	0	0	6	59	35	2	6	4	

Codes distribution



How many clusters?

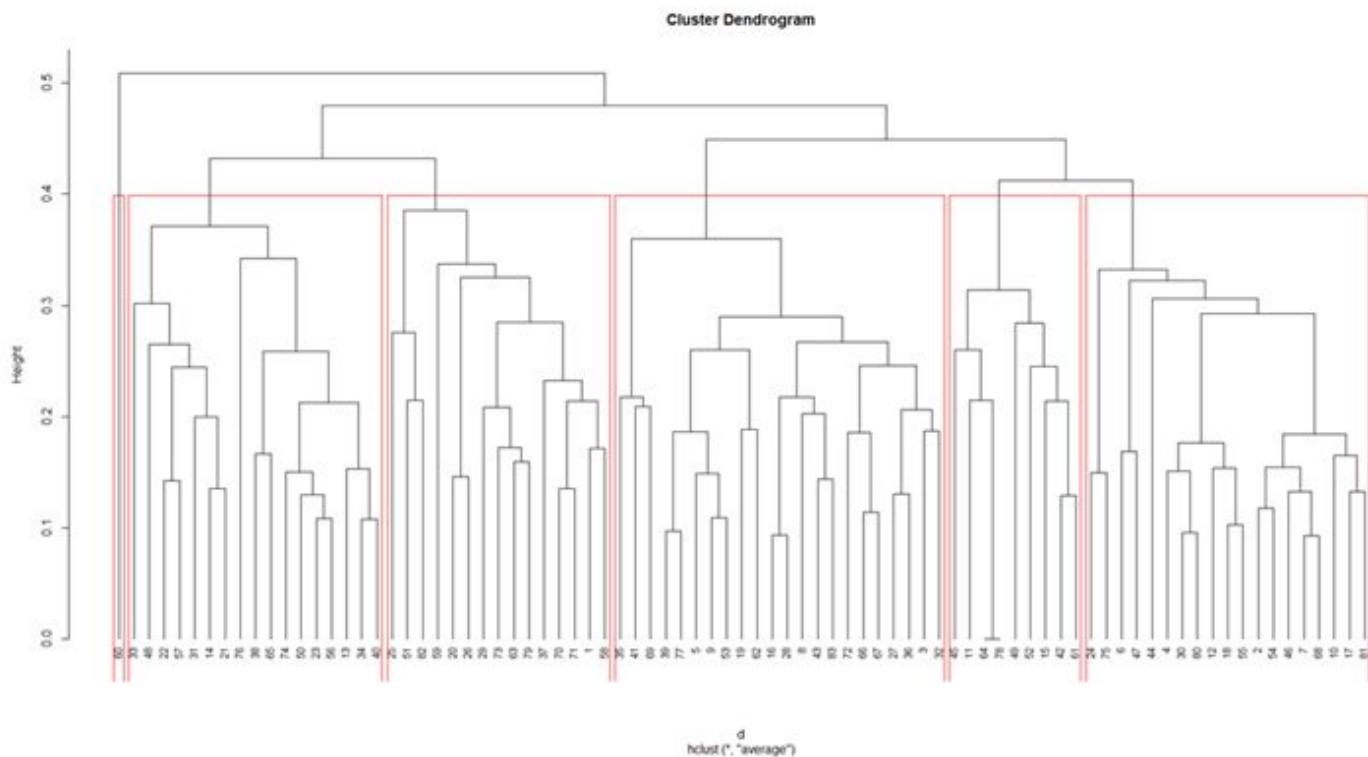
* Among all indices:

- * 3 proposed 2 as the best number of clusters
- * 4 proposed 3 as the best number of clusters
- * 1 proposed 4 as the best number of clusters
- * 1 proposed 5 as the best number of clusters
- * 6 proposed 6 as the best number of clusters
- * 1 proposed 7 as the best number of clusters
- * 3 proposed 9 as the best number of clusters
- * 1 proposed 13 as the best number of clusters
- * 3 proposed 20 as the best number of clusters

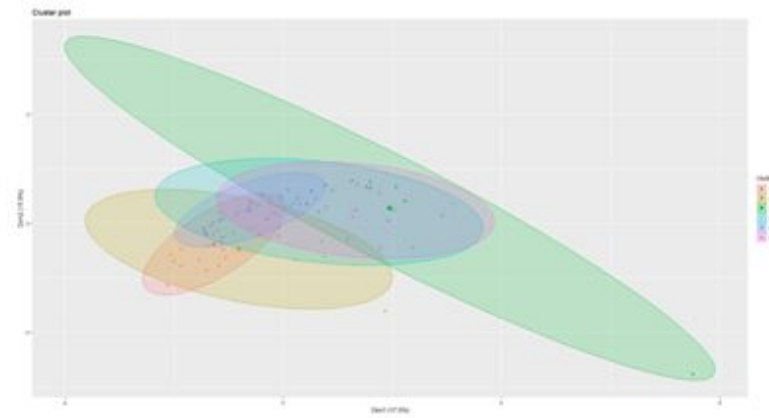
***** Conclusion *****

* According to the majority rule, the best number of clusters is 6

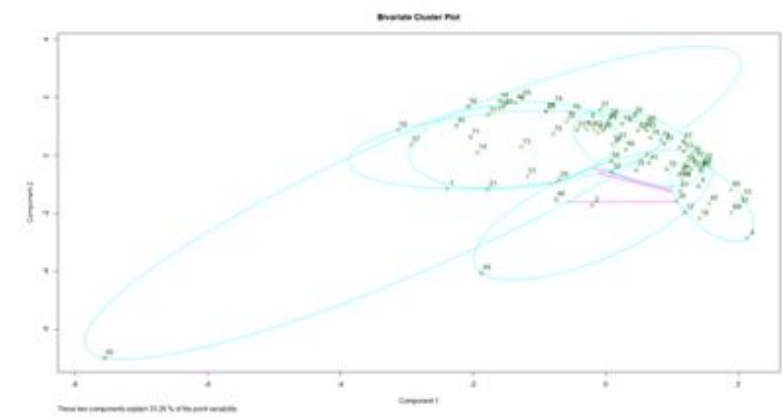
Comparing Models



1	2	3	4	5	6
15	19	22	9	17	1



1	2	3	4	5	6
14	16	10	13	19	11

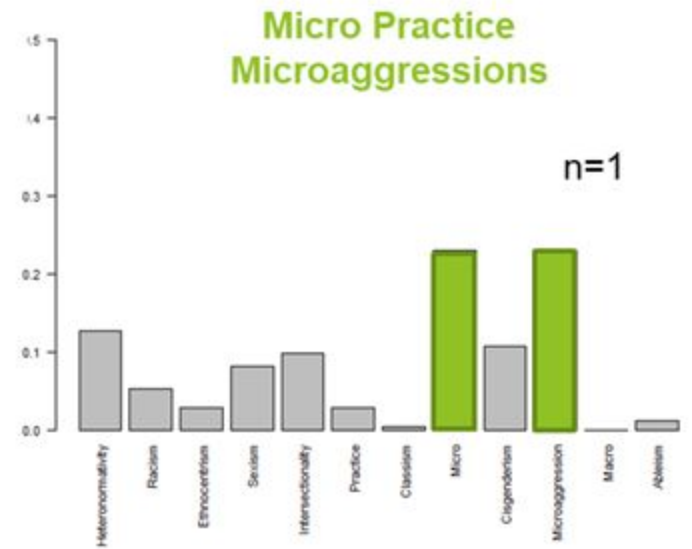
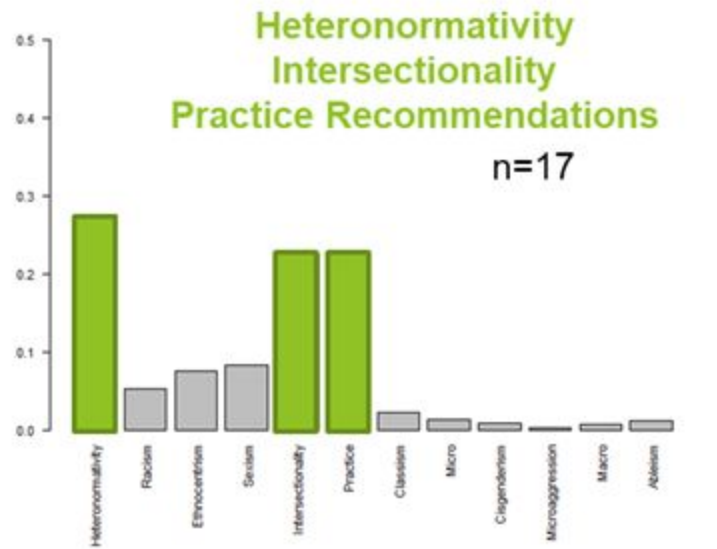
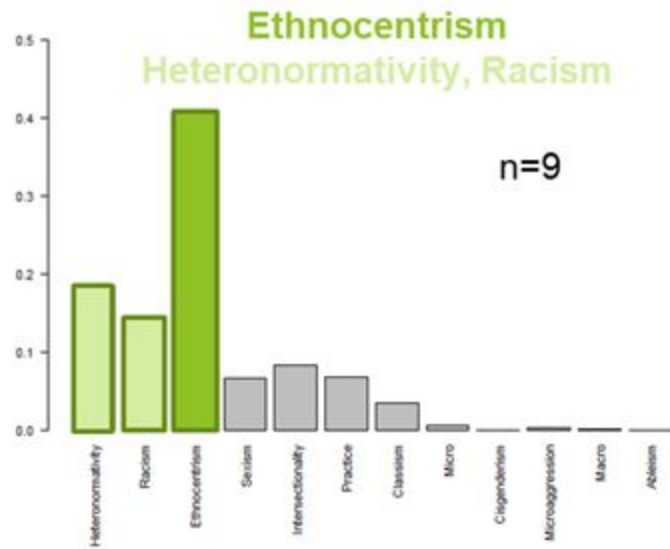
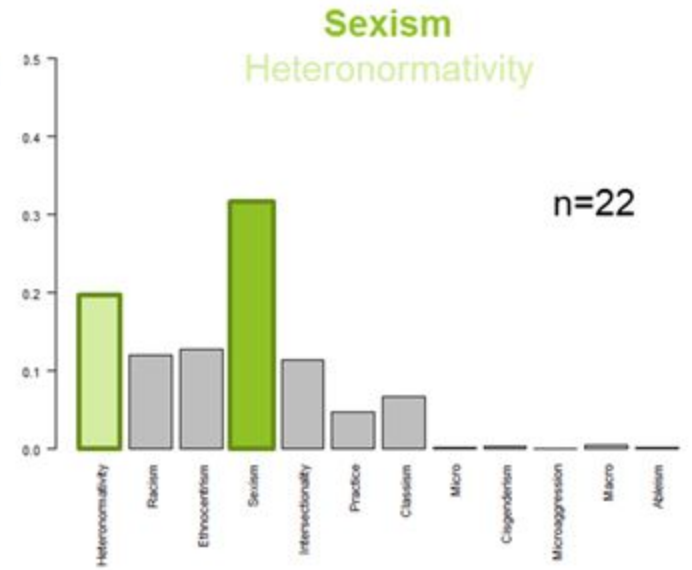
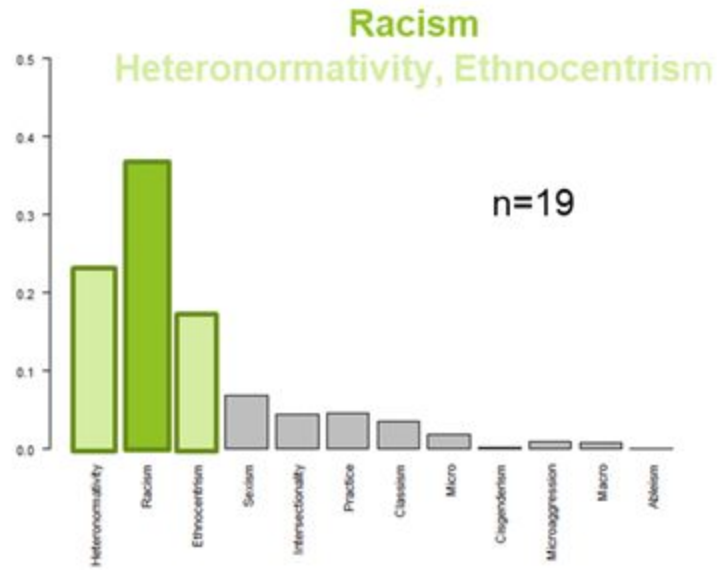
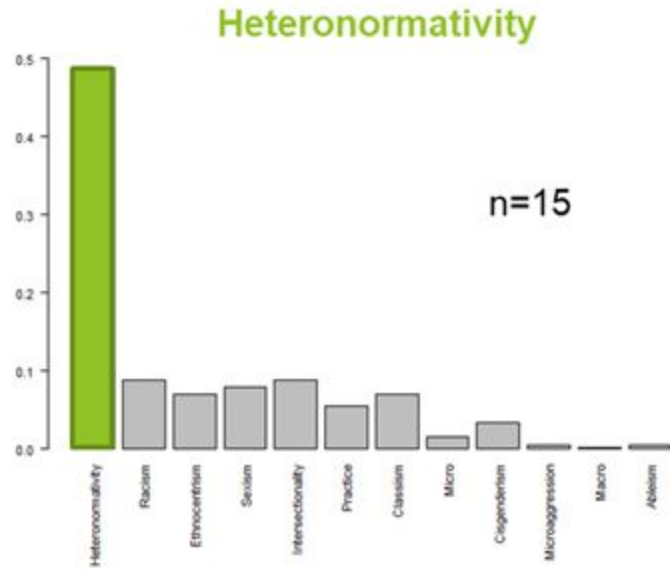


1	2	3	4	5	6
16	17	15	13	8	14

Confirming Model Choice

Codes	Hierarchical Model p	KMeans Model p	PAM Model p
Intersectionality	<0.001	<0.001	<0.001
Macro	0.374	0.870	0.397
Micro	<0.001	0.346	0.692
Microaggression	<0.001	0.430	0.605
Ableism	0.007	0.002	0.020
Cisgenderism	0.003	0.085	0.106
Ethnocentrism	<0.001	<0.001	<0.001
Heteronormativity	<0.001	<0.001	<0.001
Classism	0.187	0.103	0.277
Sexism	<0.001	<0.001	<0.001
Racism	<0.001	<0.001	<0.001
Practice	<0.001	<0.001	<0.001

Hierarchical Clusters



Hierarchical Clusters: Sample Articles

1.
Heteronormativity

Exploring Biculturality and Beauty Standards in Sexual Minority Women, *Psychology of Sexual Orientation and Gender Diversity* (2015)

2.
Racism
Heteronormativity, Ethnocentrism

"Antiblackness in Education Policy and Discourse", *Theory Into Practice* (2016)

3.
Sexism
Heteronormativity

"Capitalism and Welfare Reform: Who Really Benefits?", *Race, Gender and Class* (2008)

4.
Ethnocentrism
Heteronormativity, Racism

"Discrimination and Substance Use Disorders among Latinos", *American Journal of Public Health* (2014)

5.
Heteronormativity
Intersectionality
Practice Recommendations

Intersectionality and Social Work: Omissions of Race, Class, and Sexuality in Graduate School Education, *Journal of Social Work Education* (2016)

6.
Micro Practice
Microaggressions

Sexual Orientation and Gender Identity Microaggression, *Journal of Ethnic and Cultural Diversity in Social Work* (2017)

Questions / Discussion