

# Latent Factors and Dimensionality Reduction

INST414 - Data Science Techniques

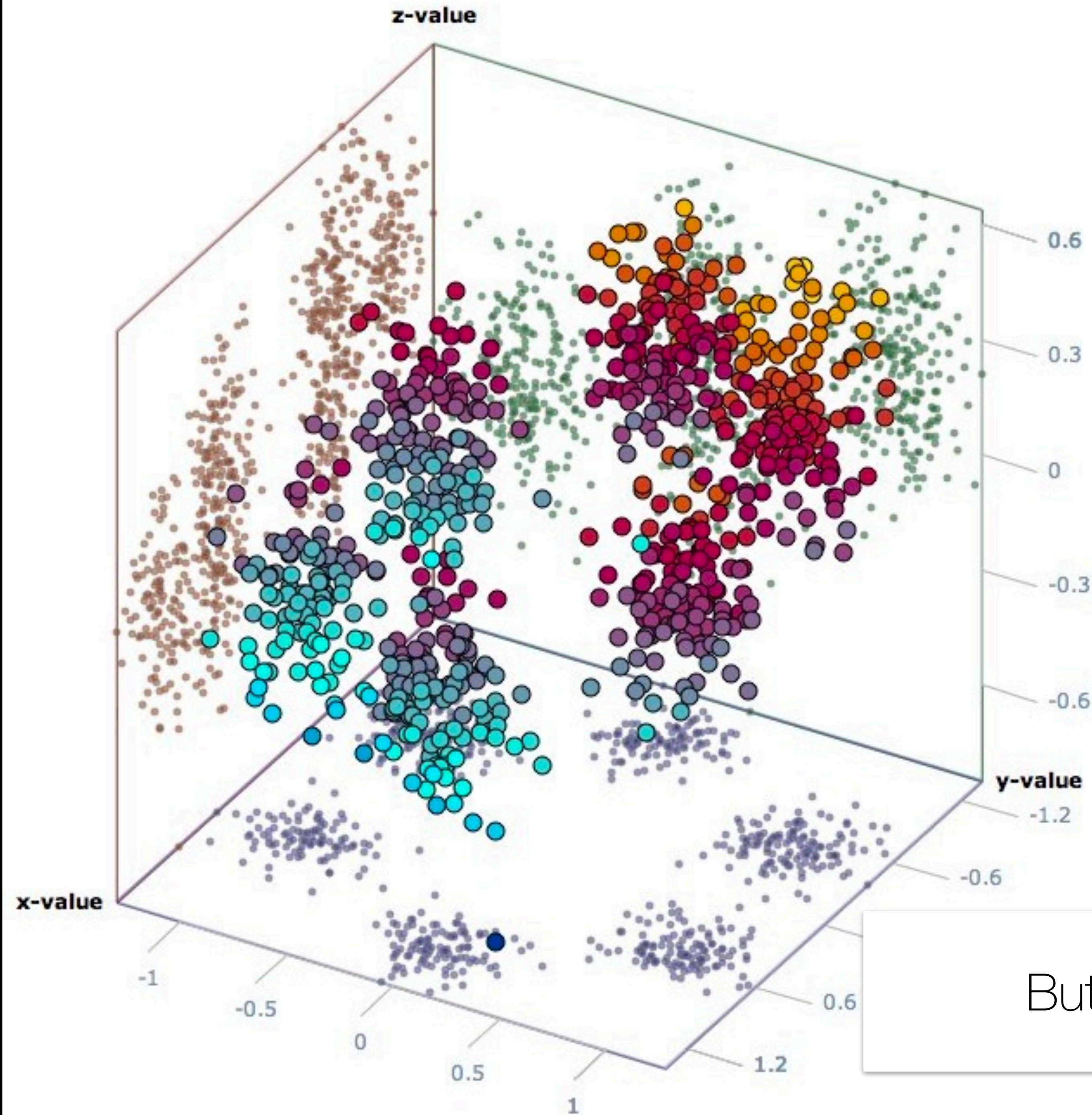
IMDb



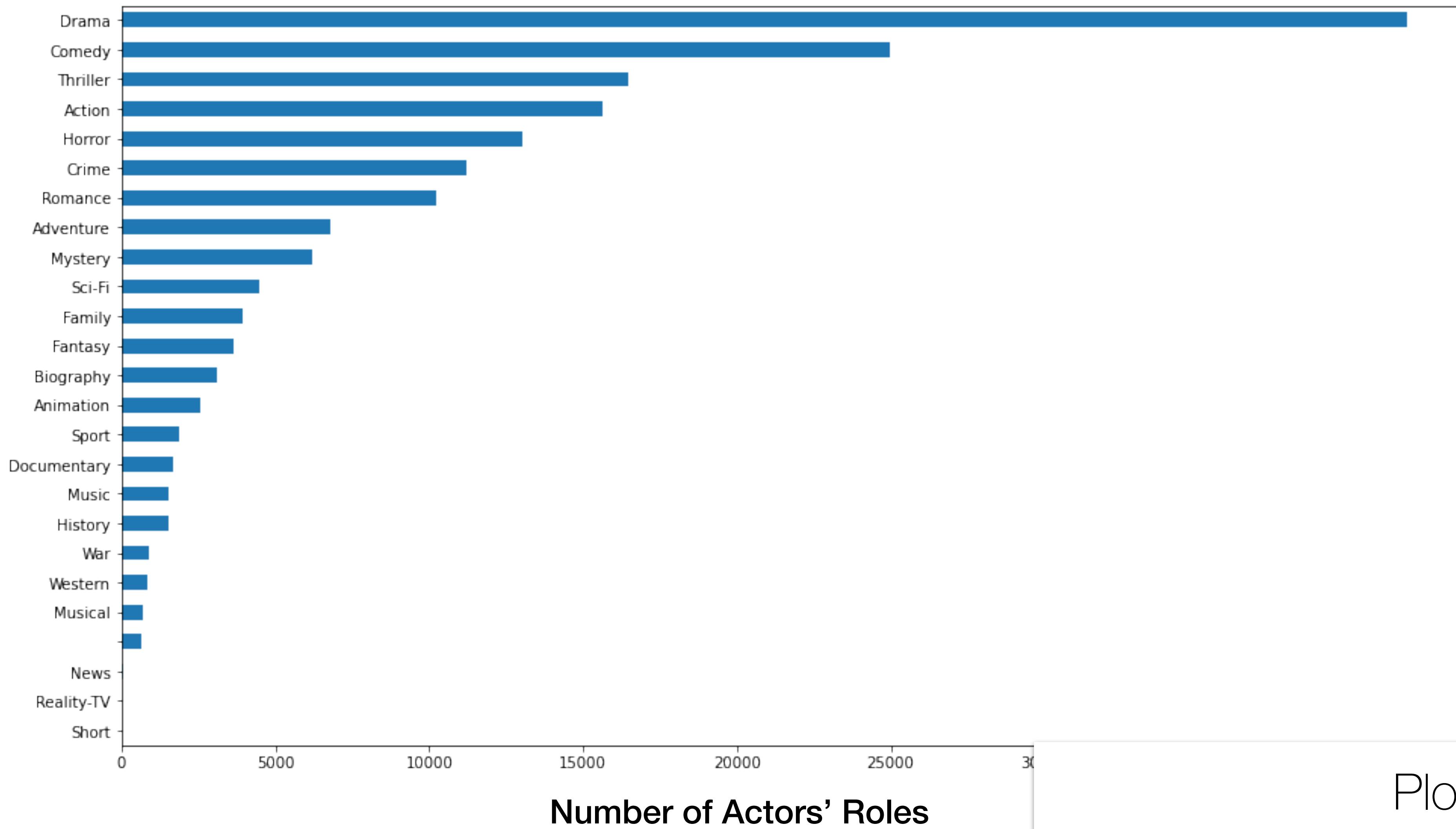
actor\_genre\_df.head(20)

	Comedy	Fantasy	Romance	Action	Crime	Adventure	Mystery	Thriller	Drama	Biography	...	Sport	News	Family	Western	Short
nm0000212	16.0	3.0	16.0	5.0	4.0	2.0	5.0	3.0	16.0	2.0	...	0.0	0.0	0.0	0.0	0.0
nm0413168	8.0	3.0	6.0	14.0	6.0	11.0	5.0	2.0	13.0	5.0	...	0.0	0.0	0.0	0.0	0.0
nm0000630	10.0	2.0	6.0	4.0	1.0	2.0	2.0	4.0	17.0	6.0	...	4.0	1.0	1.0	0.0	0.0
nm0005227	12.0	1.0	3.0	2.0	0.0	3.0	0.0	1.0	5.0	1.0	...	1.0	0.0	0.0	0.0	0.0
nm0697338	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm1300519	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0940707	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0625977	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0792032	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0496571	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm2868805	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm2866192	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0001379	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2.0	1.0	...	1.0	0.0	0.0	1.0	0.0
nm0462648	1.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	2.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0000953	6.0	0.0	0.0	1.0	3.0	0.0	0.0	2.0	9.0	7.0	...	0.0	0.0	0.0	0.0	0.0
nm0001782	0.0	0.0	0.0	1.0	0.0	1.0	0.0	1.0	1.0	1.0	...	0.0	0.0	0.0	0.0	0.0
nm0005077	1.0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	1.0	0.0
nm0550626	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0177016	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0907480	0.0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0

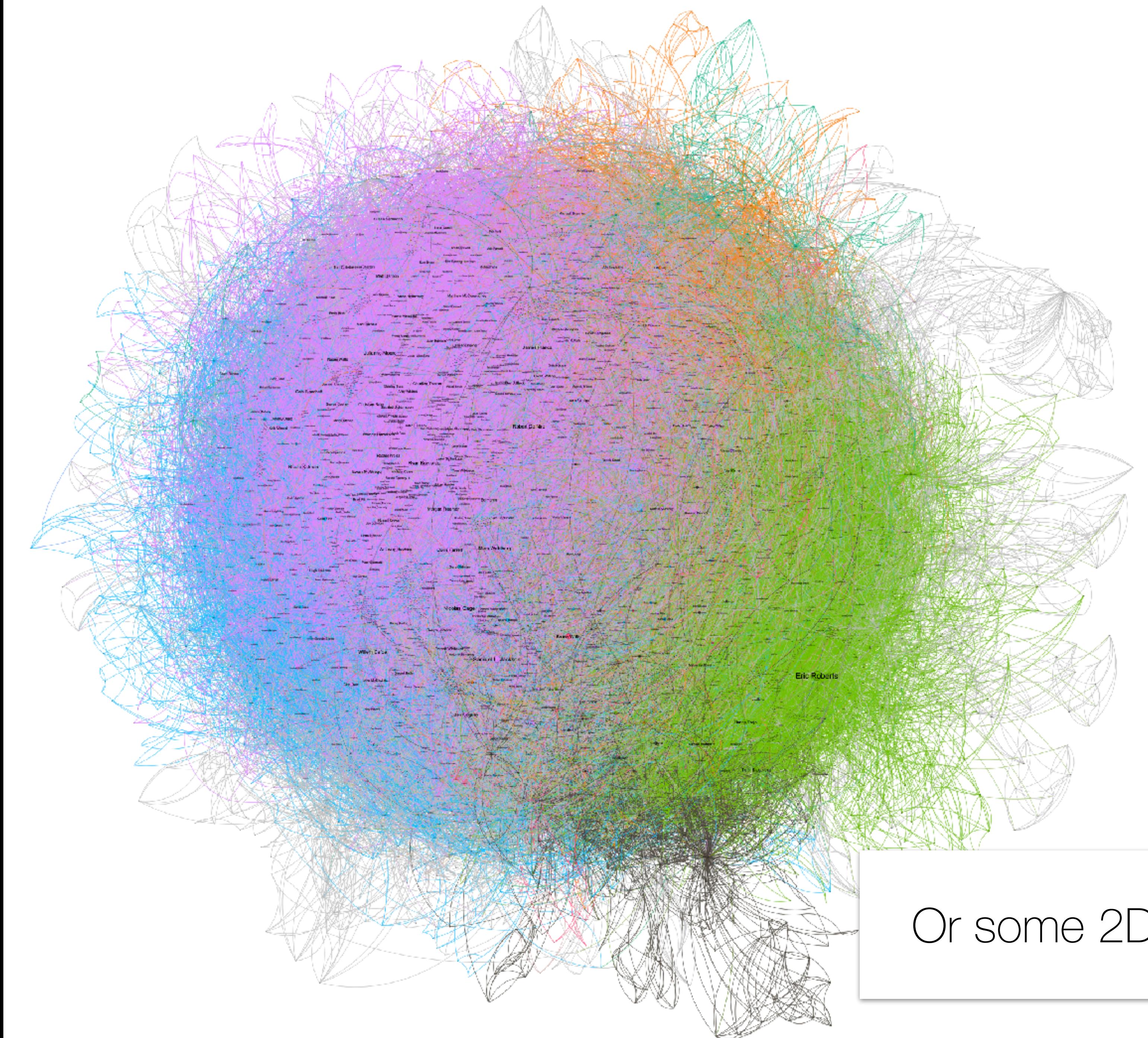
How might you visualize  
these rows?



But we have 25 genres...



Plot only the top  
two most prolific genres?



Or some 2D space of similarity

# Today's Learning Objectives

Define the connection between latent factors and dimensionality reduction

---

Describe at least two methods for dimensionality reduction

# Today's Learning Objectives

Define the connection between latent factors and dimensionality reduction

Describe at least two methods for dimensionality reduction

Core Motivation: We want to find similar content in low-dimensional spaces

---

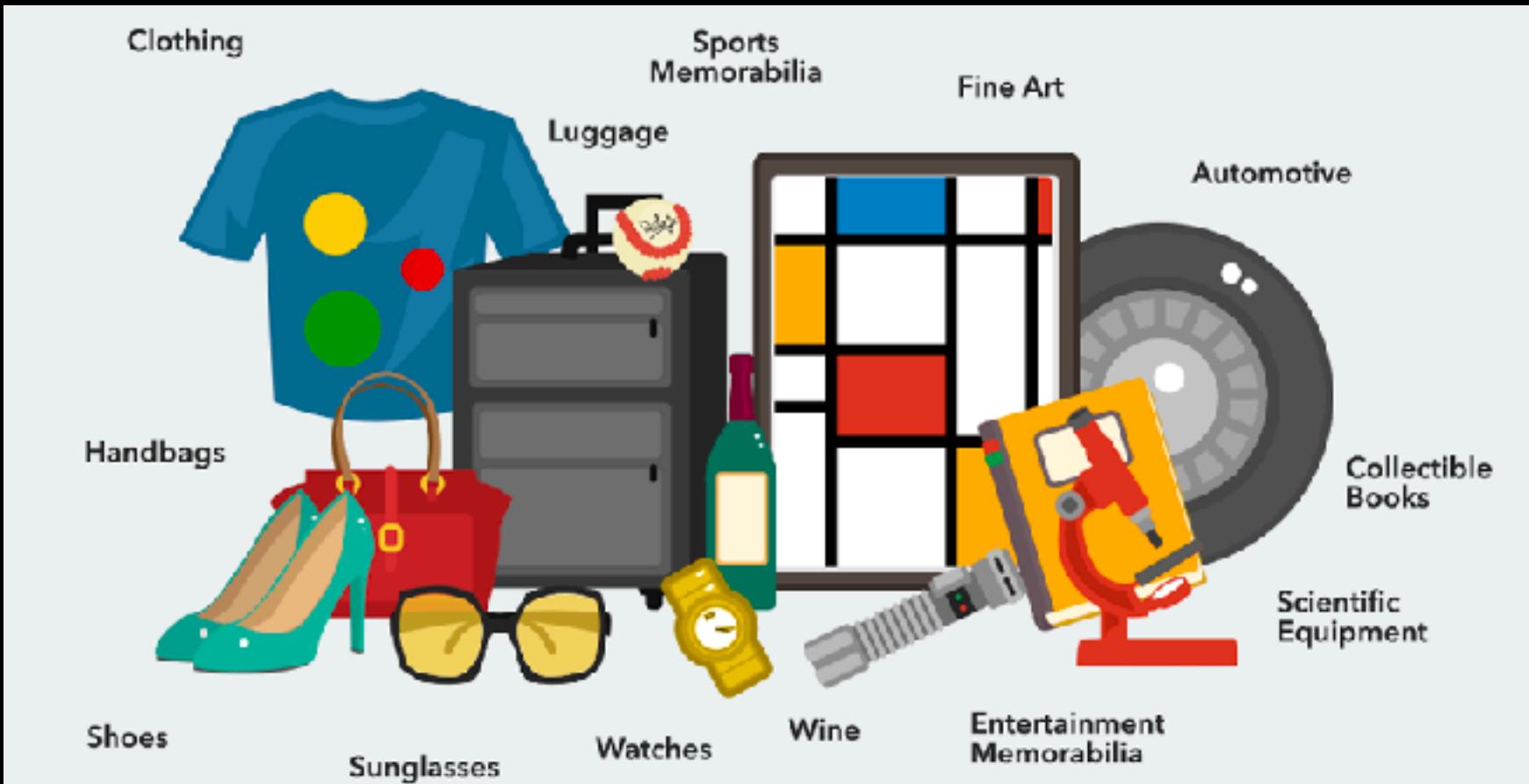
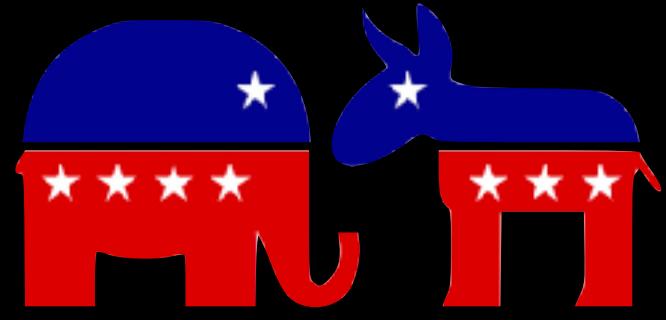
One solution is to group elements or dimensions

---

What do these “groups” mean?

## Clustering:

Identify similar groups, where  
*similar* has a real-world meaning  
( $k$  is not pre-specified)



Groups are “similar” things

# Retweets

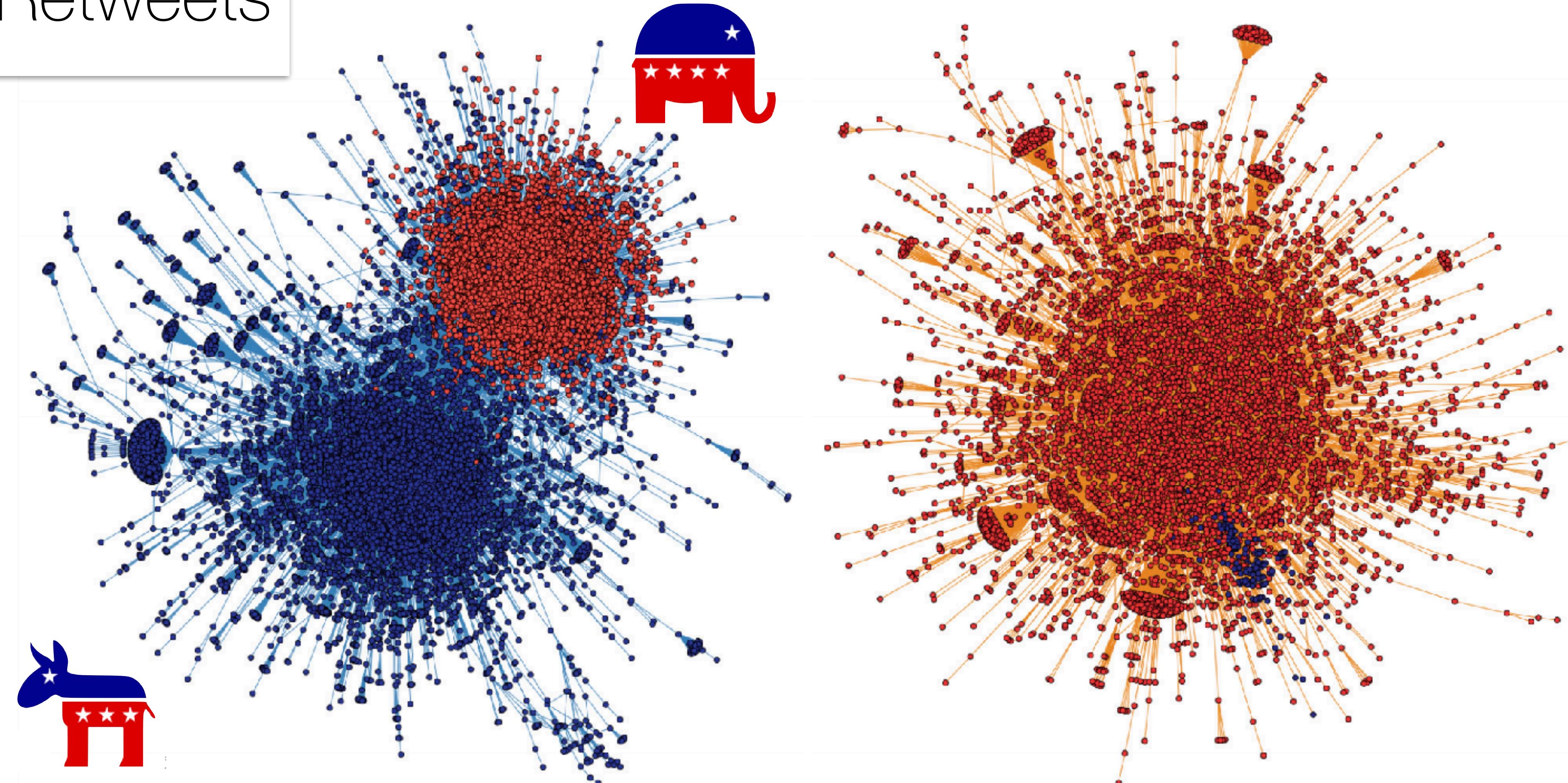


Figure 1: The political retweet (left) and mention (right) networks, laid out using a force-directed algorithm. Node colors reflect cluster assignments (see § 3.1). Community structure is evident in the retweet network, but less so in the mention network. We show in § 3.3 that in the retweet network, the red cluster A is made of 93% right-leaning users, while the blue cluster B is made of 80% left-leaning users.



Automotive

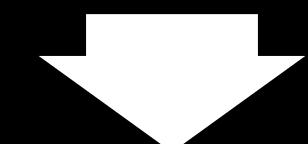
Animals

Software



**Adjacency Matrix  $A$**

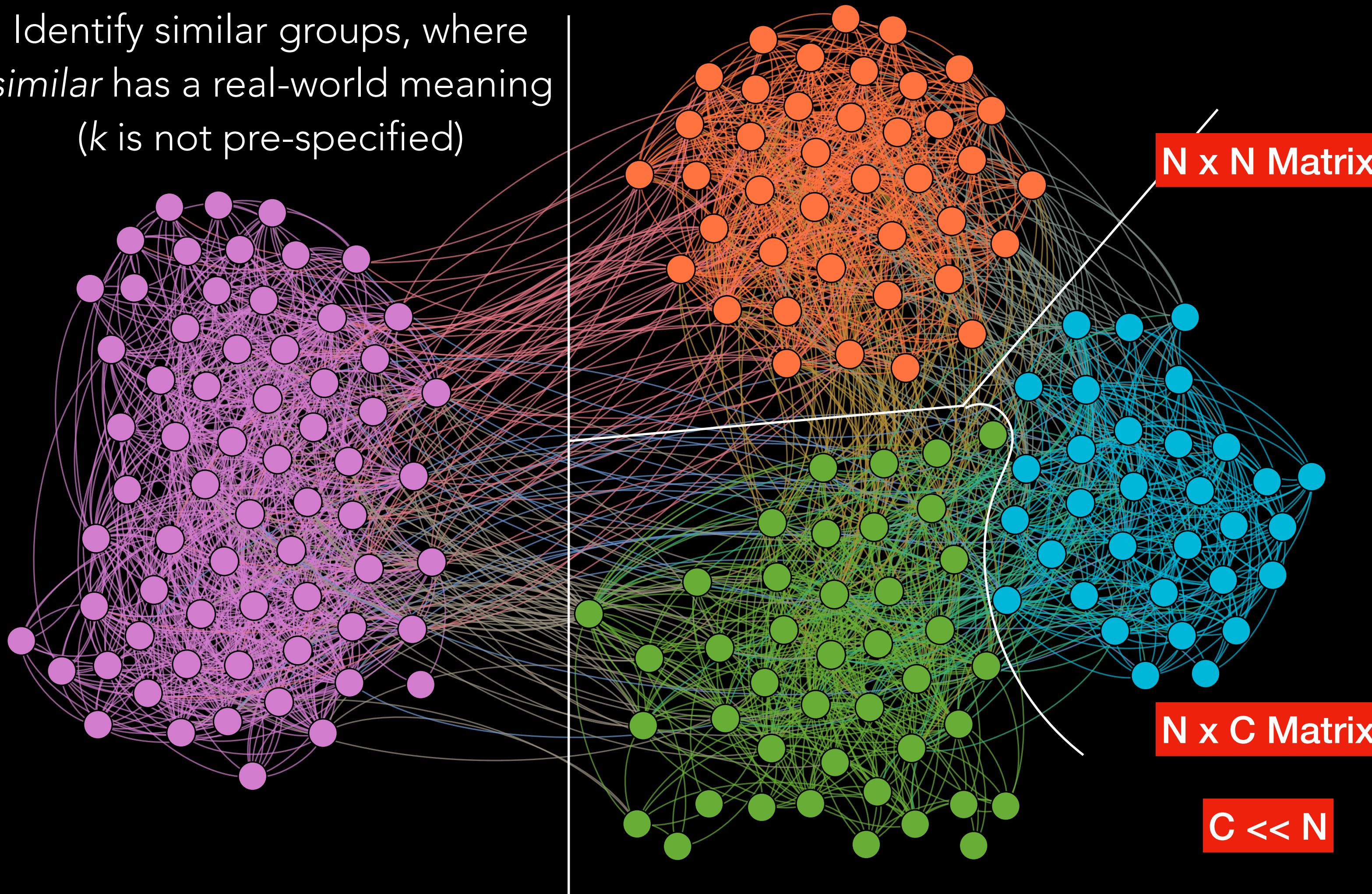
	$V_1$	$V_2$	$V_3$	$\dots$	$V_n$
$V_1$	0	0	1		1
$V_2$	1	0	0		0
$V_3$	0	1	0		0
$\dots$					
$V_n$	0	0	0		0



	$C_1$	$C_2$	$C_3$	$C_4$
$V_1$	0	0	1	0
$V_2$	1	0	0	0
$V_3$	0	1	0	0
$\dots$				
$V_n$	0	0	0	1

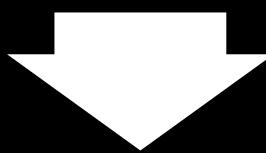
## Clustering:

Identify similar groups, where  
*similar* has a real-world meaning  
( $k$  is not pre-specified)

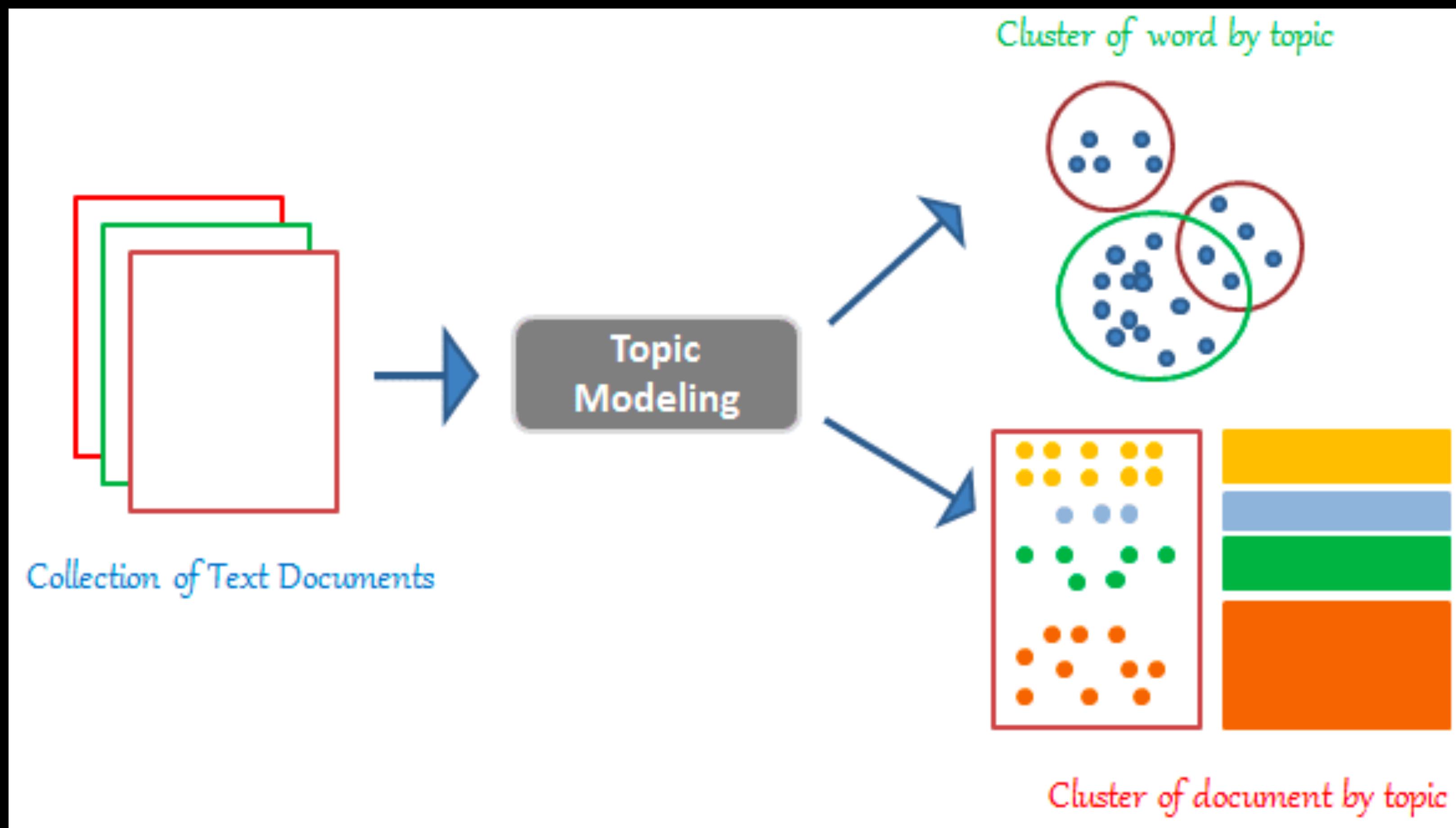


## Bag-of-Words/Shingle Matrix

	W <sub>1</sub>	W <sub>2</sub>	W <sub>3</sub>	...	W <sub>n</sub>
D <sub>1</sub>	0	0	1		1
D <sub>2</sub>	1	0	0		0
D <sub>3</sub>	0	1	0		0
...					
D <sub>n</sub>	0	0	0		0



	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>	T <sub>4</sub>
D <sub>1</sub>	0	0	1	0
D <sub>2</sub>	1	0	0	0
D <sub>3</sub>	0	1	0	0
...				
D <sub>n</sub>	0	0	0	1



Topic modeling and Latent Dirichlet Allocation

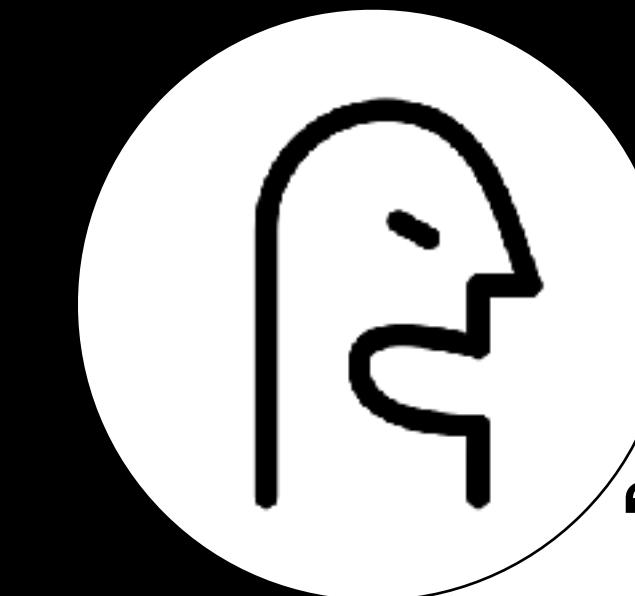
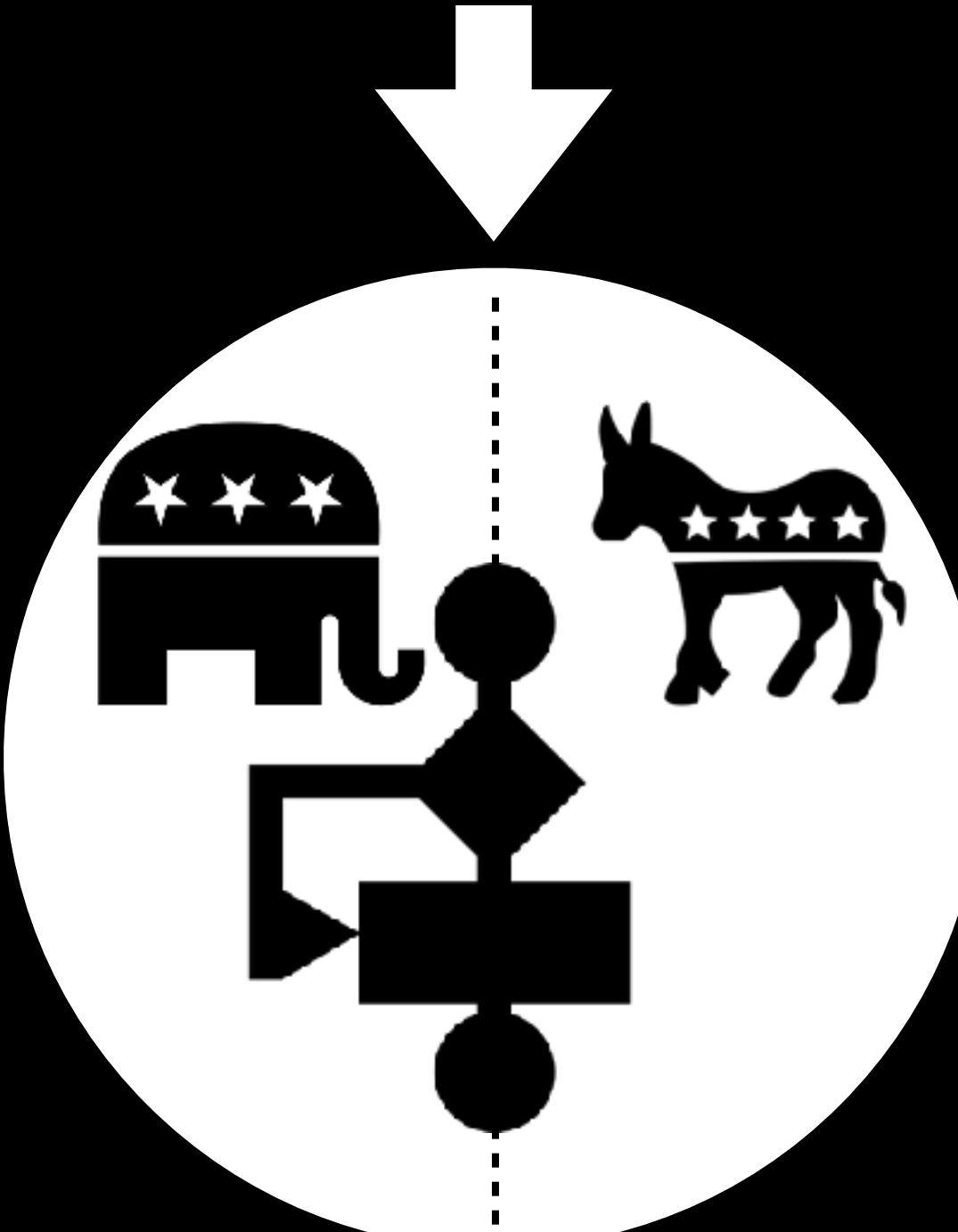
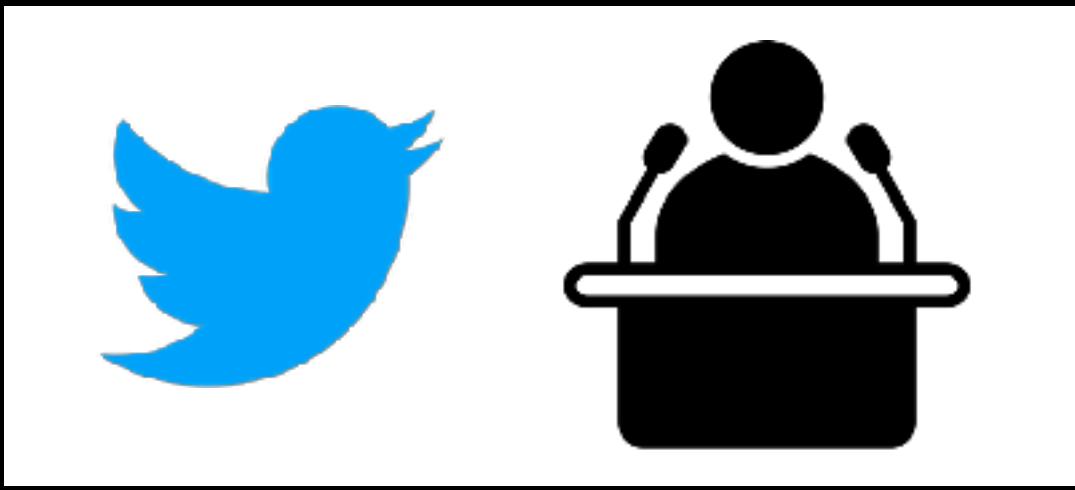
In some sense, you're familiar with dimensionality reduction already

These groups capture hidden ("latent") relationships

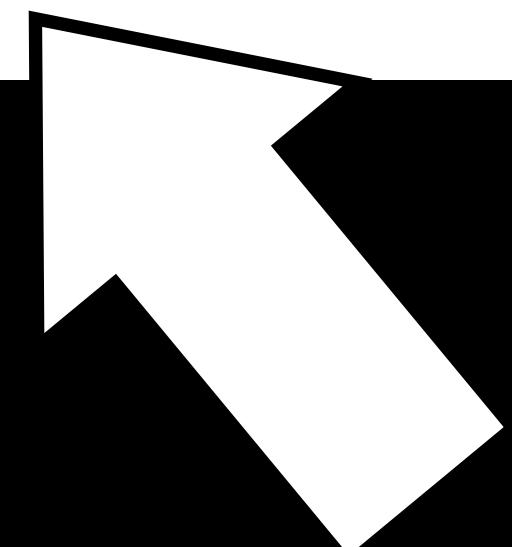
May not know what these latent factors are a priori



# Link Sharing and Ideology

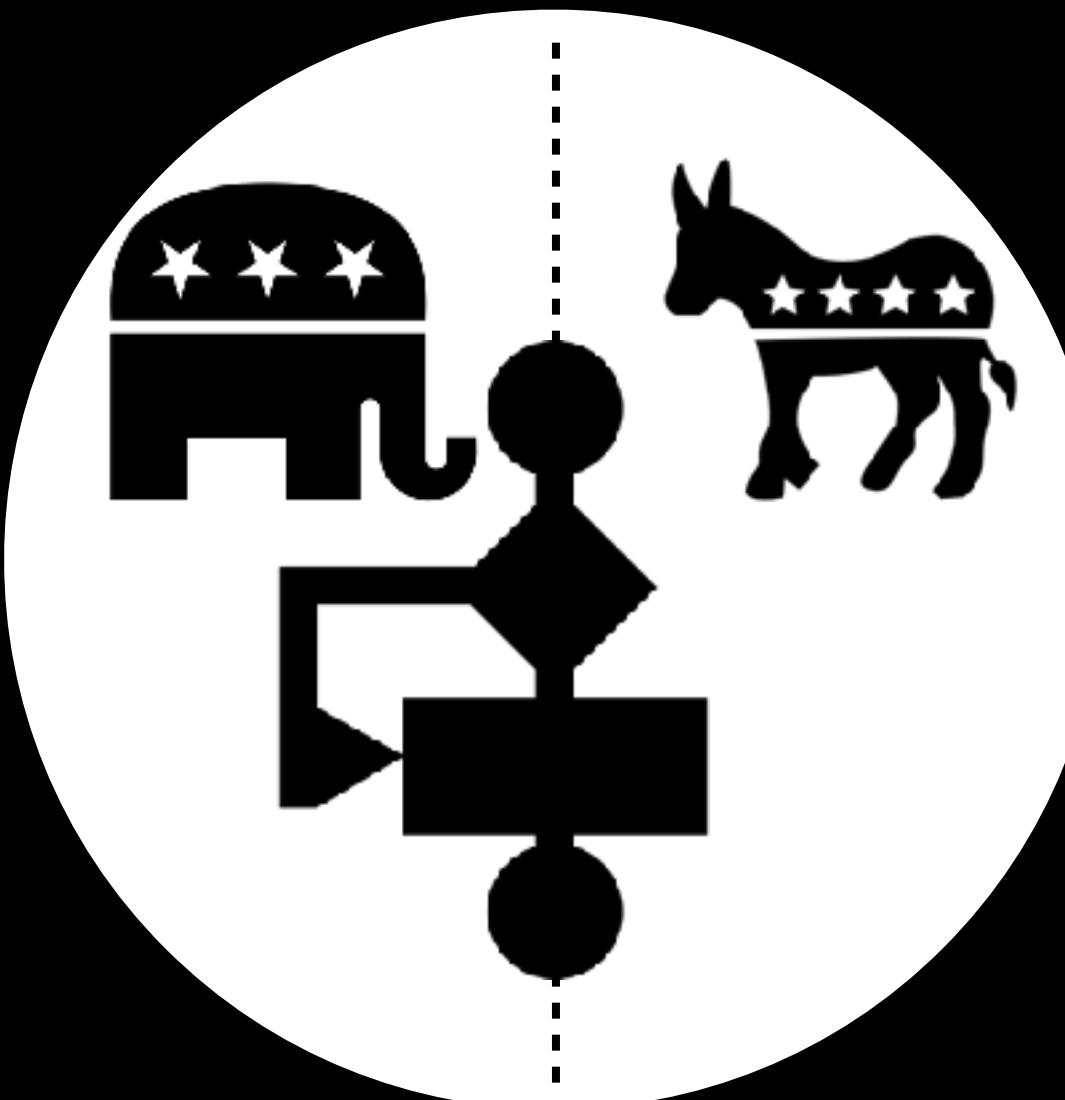
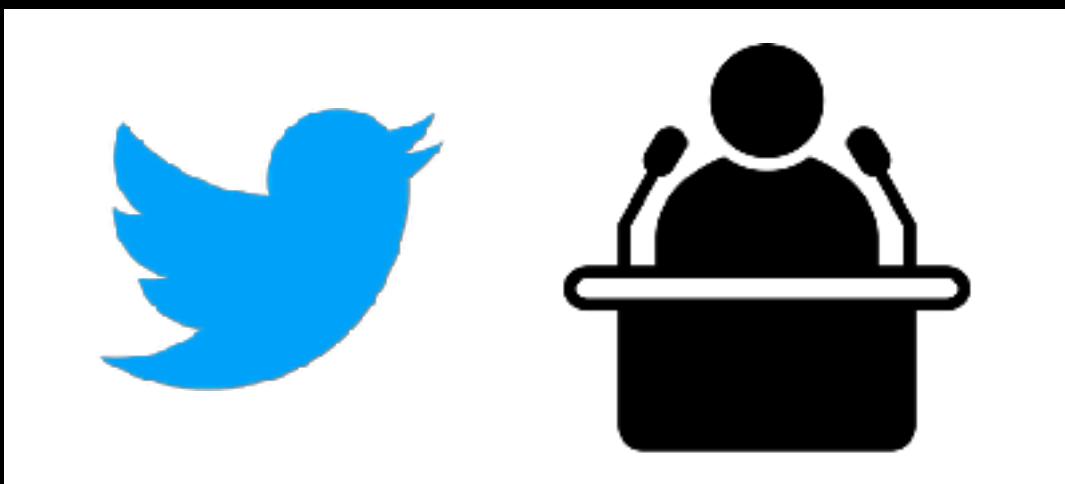


**“Amazon pulls out of New York City deal, leaving many to wonder whether this was a net positive or negative.” - [NYTimes.com](#)**





# Link Sharing and Ideology



Elizabeth Warren  @ewarren · 3h

Sen. Lisa Murkowski  @lisam  
I have long been advocating for t

Cory Gardner Retweeted

Posts

Elizabeth Warren  March 21 at 2:35  
I don't take PAC money. I take grassroots movement to take all of us fighting together.

SECURE.ACTBLUE.COM  
Donate to Warren

Lisa Murkowski  
March 15 at 12:32 PM · 

While I was at #CERAWEEK, a topic that frequently came up in conversation is addressing climate change with future energy policy. Climate change is real and we must take reasonable and practicable action to address it. But the best way to address this issue is through major technological breakthroughs. The U.S. is a global leader on technology, which is our best path forward. Whether it's increasing efficiencies or promoting advanced nuclear, next-generation energy storage, and carbon capture, these are the areas we need to focus on to combat our changing climate. We need to leave behind the sensationalism and the sloganizing and work on bipartisan genuine solutions.

END POVERTY

GREEN DEAL

CNBC.COM

Ocasio-Cortez's Green New Deal is not going over well at one of the most liberal congressional caucuses



# Link Sharing and Ideology



Elizabeth Warren



Lisa Murkowski



Bernie Sanders

*nytimes.com*  
*reuters.com*  
*bbc.co.uk*  
*vox.com*  
*HuffPost.com*  
*breitbart.com*  
...  
*foxnews.com*  
*npr.org*

23	7	9	37	22	0	...	0	20
----	---	---	----	----	---	-----	---	----

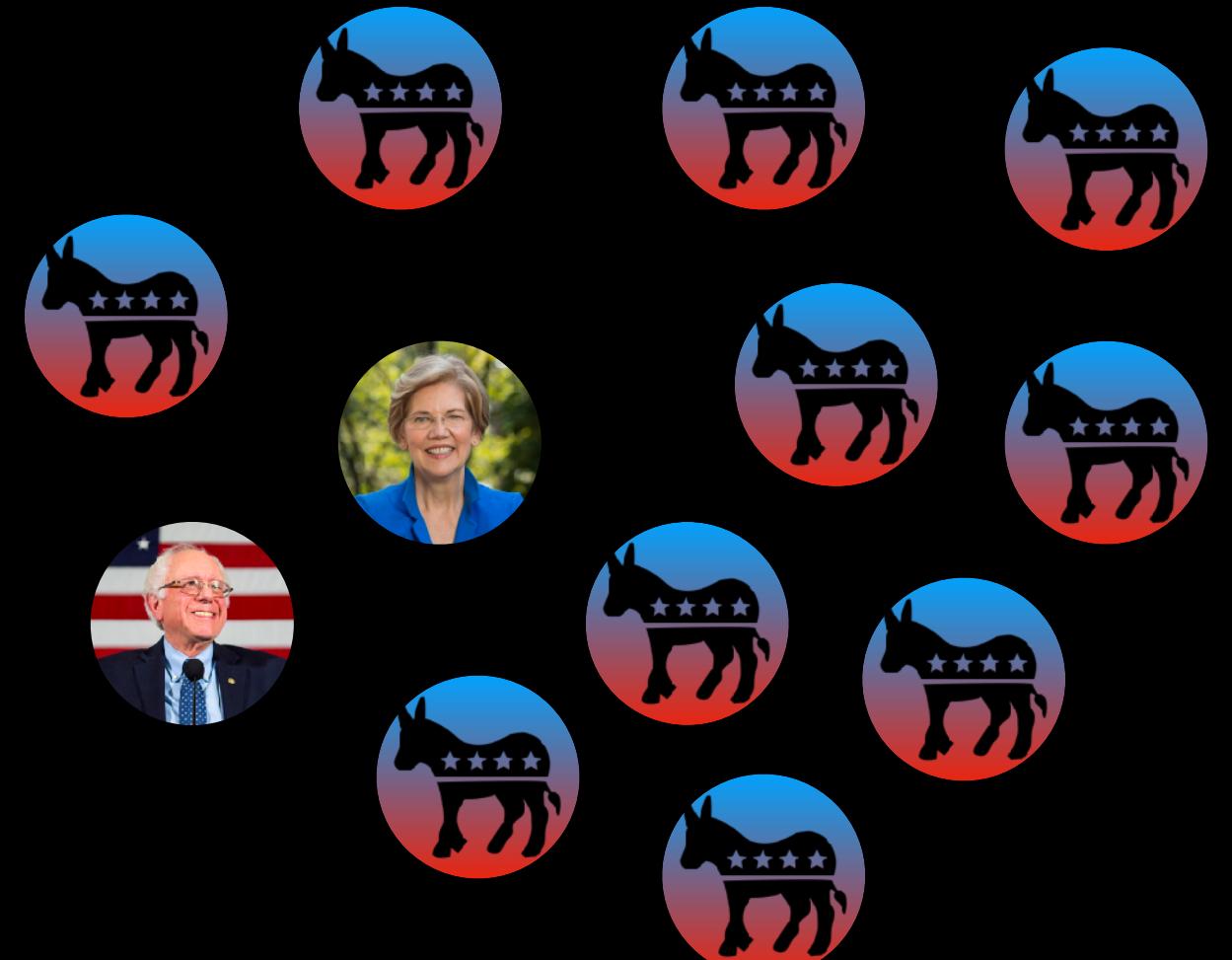
3	11	10	0	1	5	...	21	3
---	----	----	---	---	---	-----	----	---

31	23	17	40	39	0	...	0	17
----	----	----	----	----	---	-----	---	----

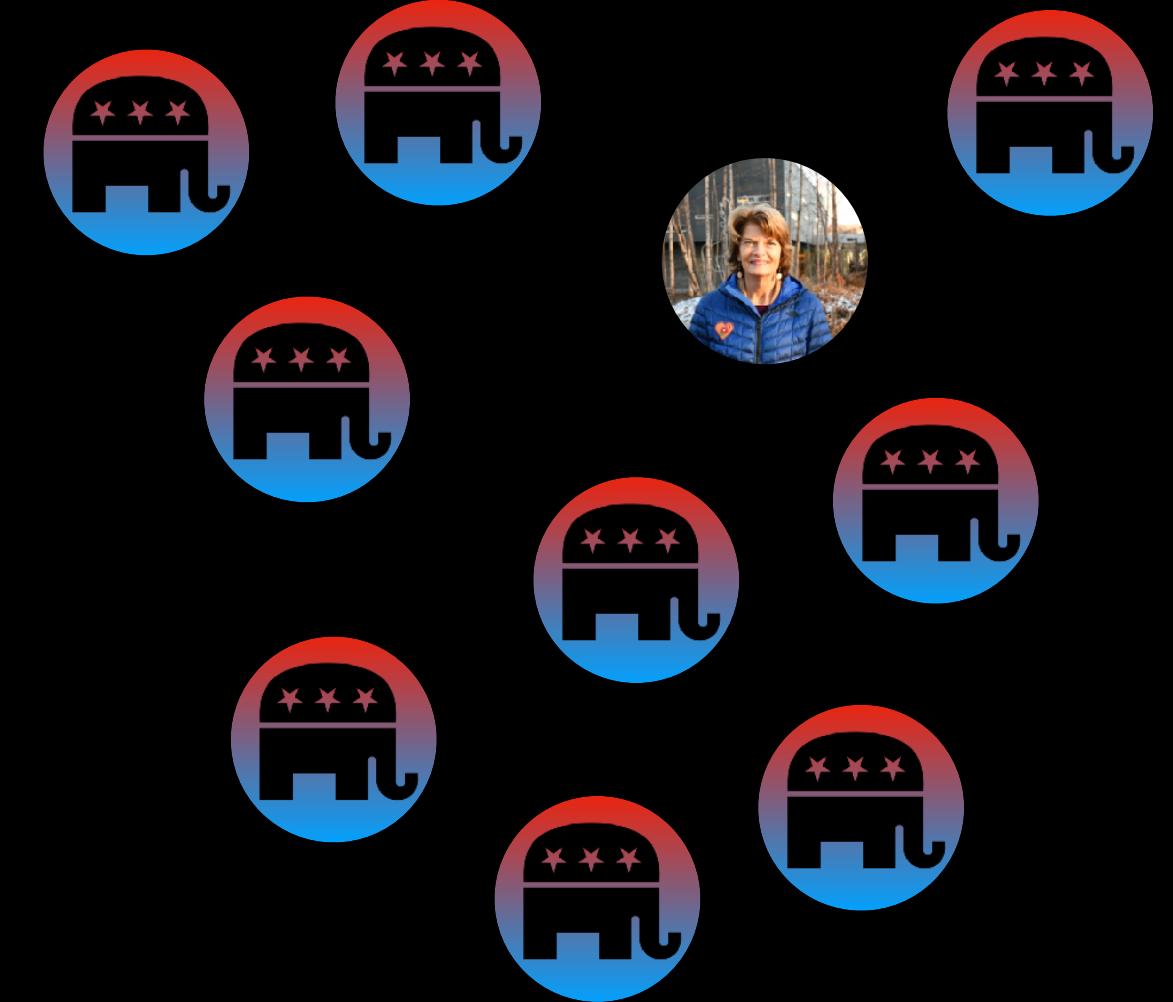


# Link Sharing and Ideology

One dimension  
for each web domain  
you share



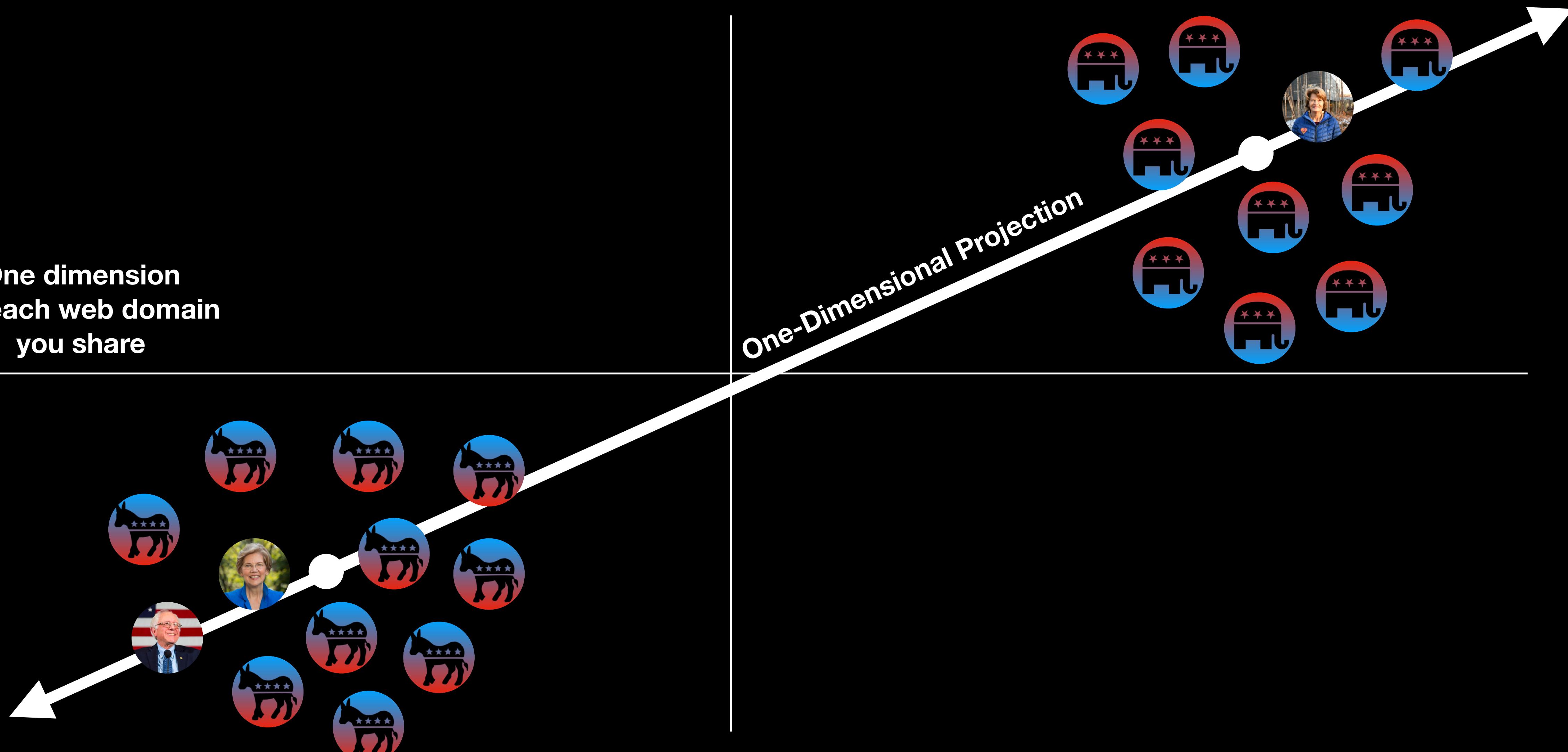
Point-clouds for various politicians' domain sharing





# Link Sharing and Ideology

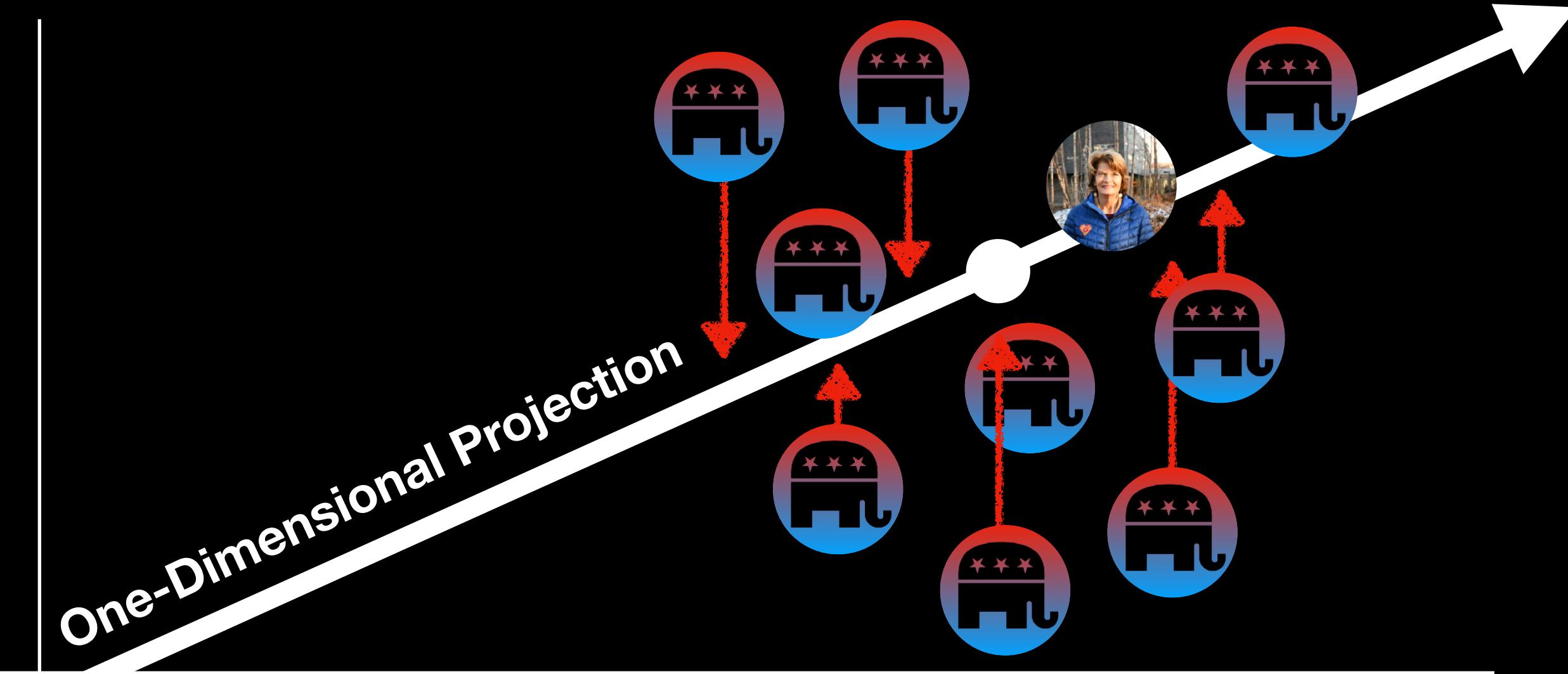
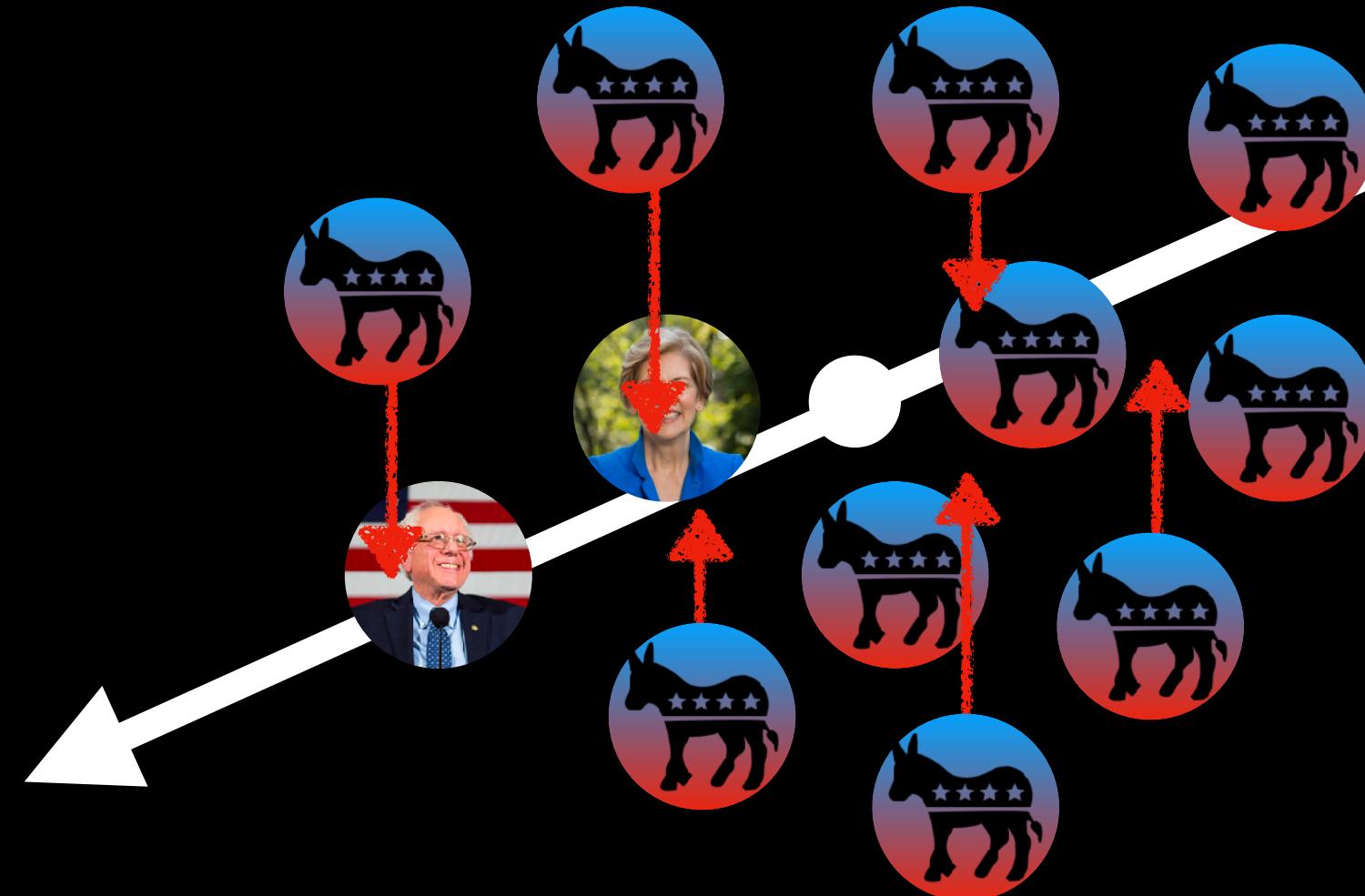
One dimension  
for each web domain  
you share





# Link Sharing and Ideology

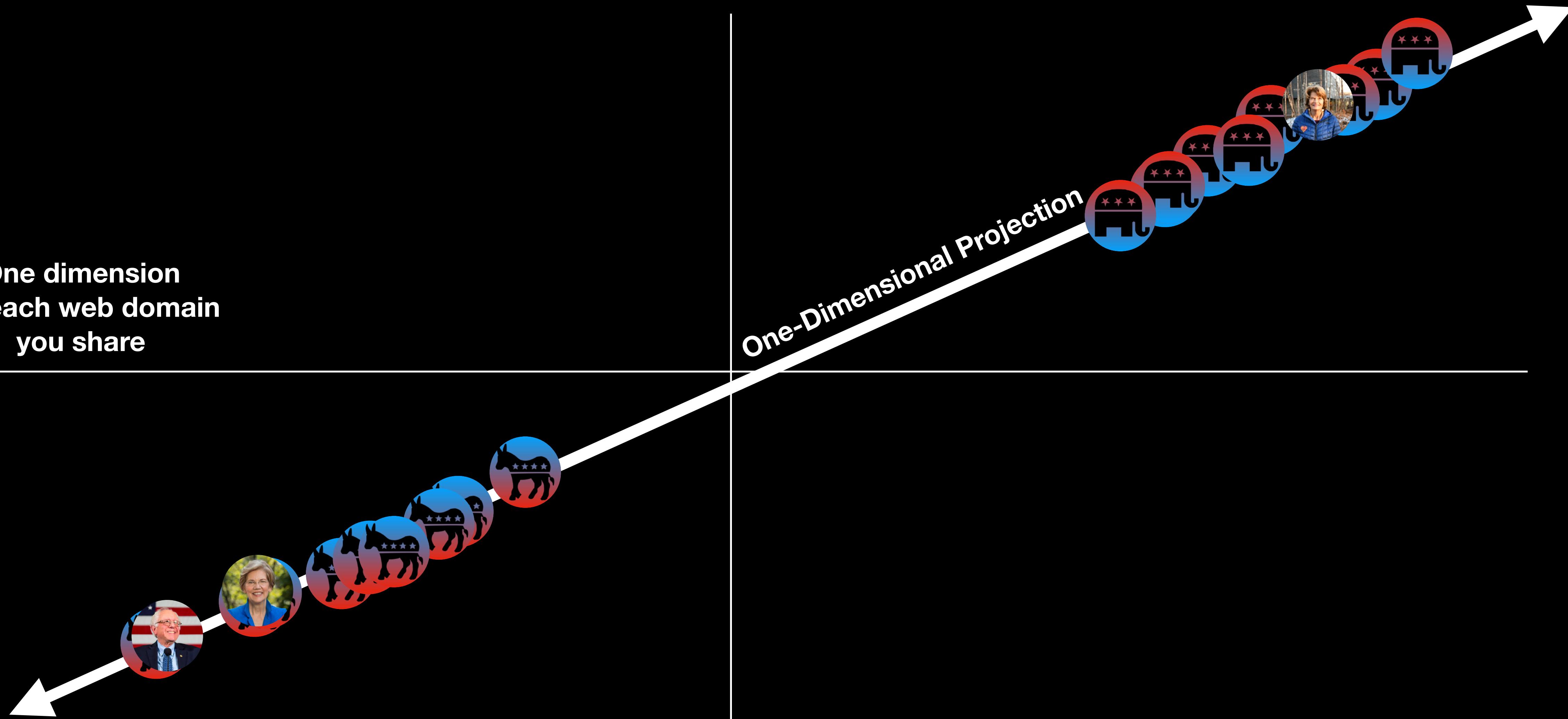
One dimension  
for each web domain  
you share





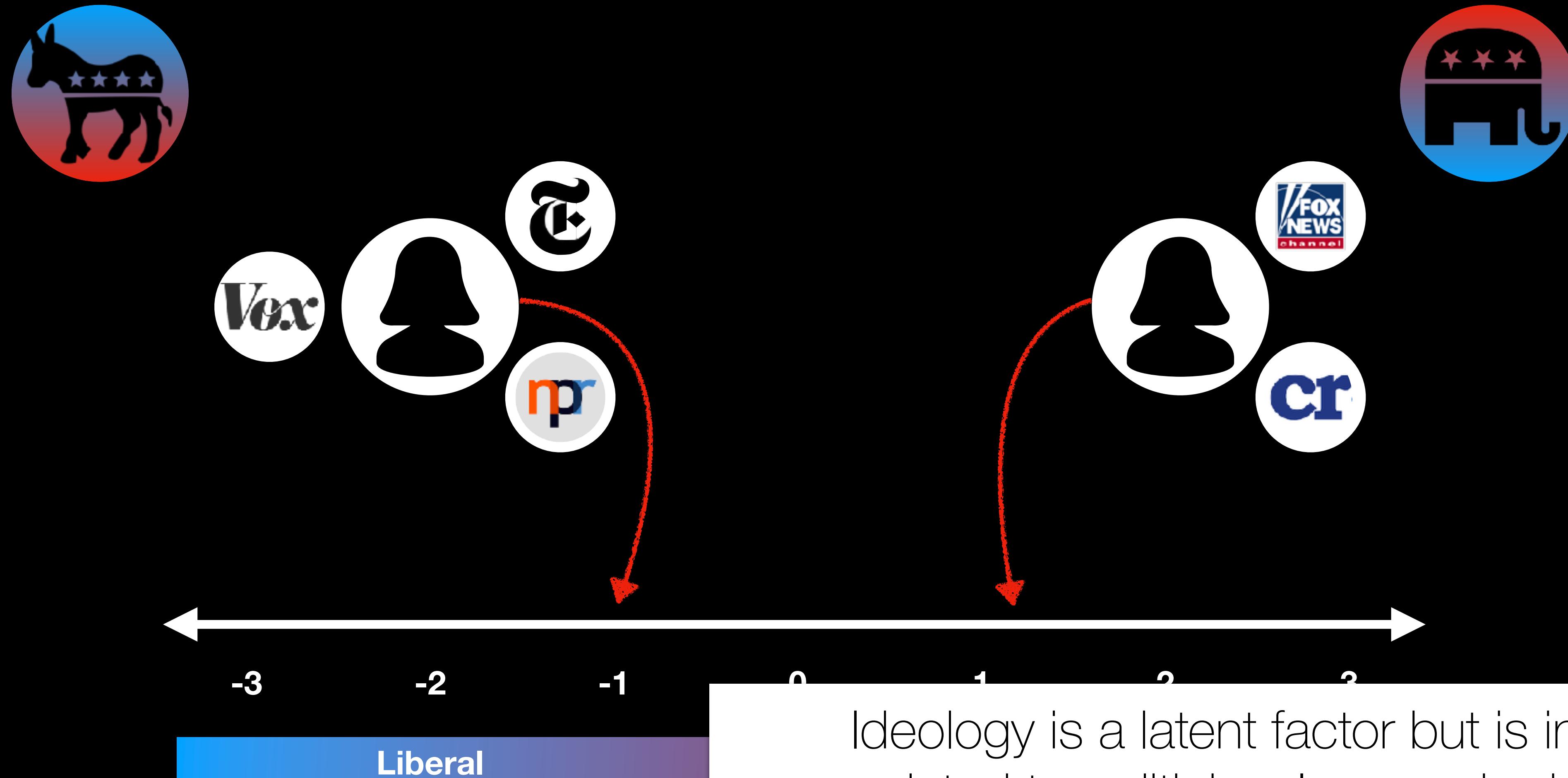
# Link Sharing and Ideology

One dimension  
for each web domain  
you share





# Link Sharing and Ideology



1. thing, bad, guy, happen, lot, worry, imagine, good, forget, shit -> (UNKNOWN)
2. love, stay at home, good, hope, new normal, morning, stay at home, night, situation, care to stay, never
3. work, great, team, hard, amazing, time, idea, part, proud, job -> Appreciation
4. health, care, patient, hospital, public, doctor, medical, system, nurse, treat -> Healthcare
5. virus, people, die, corona, kill, infect, person, thousand, black, dead -> Deaths during covid
6. case, death, confirm, report, number, total, rise, record, recover, late -> Number of cases/deaths
7. week, lockdown, start, end, level, move, restriction, daily, covid, begin -> Extension of lockdown
8. mask, face, wear, hand, wash, water, touch, cover, clean, glove -> Precautions
9. bus, today, join, cancer, live, pm, tonight, tomorrow, event, covid -> Cancellation of events
10. time, change, learn, human, important, experience, manage, create, covid, tip -> Survival tips
11. pay, money, job, company, lose, cut, month, employee, leave, cost -> Layoffs/Pay-cut
12. covid, call, state, pandemic, medium, situation, social, lead, current, issue ->
13. watch, live, play, video, talk, game, start, season, set, win -> Sports on TV
14. time, year, life, back, long, month, save, wait, normal, spend -> Wait to go back to normal
15. give, people, agree, police, reason, fine, chance, power, mind, understand -> (UNKNOWN)
16. people, covid, man, woman, group, chinese, warn, affect, issue, concern -> (UNKNOWN)
17. trump, lie, wrong, fail, blame, stupid, stop, truth, vote, dangerous -> Blaming Trump
18. covid, today, feel, day, sign, good, bit, call, hope, finally -> Family Occasions
19. lockdown, day, isolation, quarantine, music, busy, happy, coronalockdown, single, dog -> Schedule
20. people, put, follow, place, rule, law, break, lock, continue, avoid -> (Language Barrier)
21. covid, share, listen, post, fact, important, message, check, explain, list -> Sharing opinions
22. government, response, plan, measure, public, action, canadian, official, announce, release -> Govt. response
23. read, hear, story, word, write, sad, thought, book, article, full -> Reading books
24. buy, big, hit, market, panic, fear, store, demand, sell, shop -> Demand and supply
25. home, stay, family, safe, friend, send, healthy, strong, distance, hope -> Safety Messages
26. news, question, travel, quarantine, return, answer, flight, citizen, break, leave -> Travel restrictions
27. make, good, thing, point, happen, decision, big, easy, light, sense -> (UNKNOWN)
28. support, business, community, local, small, provide, donate, fund, service, program -> Support local
29. treatment, vaccine, find, interesting, base, study, research, cure, scientist, science -> Vaccine
30. pandemic, crisis, global, economy, impact, economic, future, opportunity, industry, threat -> Economic
31. spread, stop, child, school, close, kid, student, parent, risk, slow -> Closing of Schools In UK
32. run, walk, city, hour, open, exercise, drive, sit, close, space -> Transportation/ Exercise
33. food, order, deliver, eat, delivery, drink, pick, animal, local, restaurant -> Food
34. test, symptom, positive, show, testing, covid, result, contact, cough, flu -> Test results/ Symptoms
35. covid, high, risk, death, infection, number, rate, low, population, datum ->
36. free, update, late, information, link, check, visit, resource, advice, online -> Information links
37. worker, fight, covid, staff, continue, protect, essential, resident, service, key -> Medical staff appreciate
38. world, country, pandemic, leader, leadership, problem, real, show, poor, rest -> Leadership opinion

Topic models puts “words” in groups

These groups may not always make sense

1. thing, bad, guy, happen, lot, worry, imagine, good, forget, shit -> (UNKNOWN)  
2. love, staysafe, stayhome, good, hope, flattenthecurve, morning, stayathome, night, beautiful-> Safe to stay home.  
3. work, great, team, hard, amazing, time, idea, part, proud, job-> Appreciation  
4. health, care, patient, hospital, public, doctor, medical, system, nurse, treat-> Healthcare  
5. virus, people, die, corona, kill, infect, person, thousand, black, dead-> Deaths during covid  
6. case, death, confirm, report, number, total, rise, record, recover, late-> Number of cases/deaths  
7. week, lockdown, start, end, level, move, restriction, early, covid, begin->Extension of lockdown  
8. mask, face, wear, hand, wash, water, touch, cover, clean, glove->Precautions  
9. due, today, join, cancel, live, pm, tonight, tomorrow, event, covid-> Cancellation of events  
10. time, change, learn, human, important, experience, manage, create, covid, tip->Survival tips  
11. pay, money, job, company, lose, cut, month, employee, leave, cost-> Layoffs/Pay-cut  
12. covid, call, state, pandemic, medium, situation, social, lead, current, issue->  
13. watch, live, play, video, talk, game, start, season, set, win->Sports on TV  
14. time, year, life, back, long, month, save, wait, normal, spend-> Wait to go back to normal  
15. give, people, agree, police, reason, fine, chance, power, mind, understand->(UNKNOWN)  
16. people, covid, man, woman, group, chinese, warn, affect, issue, concern->(UNKNOWN)  
17. trump, lie, wrong, fail, blame, s  
18. covid, today, feel, day, sign, go  
19. lockdown, day, isolation, quara  
20. people, put, follow, place, rule,  
21. covid, share, listen, post, fact, i  
22. government, response, plan, m  
23. read, hear, story, word, write, sad, thought, book, article, full-> Reading books  
24. buy, big, hit, market, panic, fear, store, demand, sell, shop->Demand and supply  
25. home, stay, family, safe, friend, send, healthy, strong, distance, hope-> Safety Messages  
26. news, question, travel, quarantine, return, answer, flight, citizen, break, leave-> Travel restrictions  
27. make, good, thing, point, happen, decision, big, easy, light, sense->(UNKNOWN)  
28. support, business, community, local, small, provide, donate, fund, service, program-> Support local businesses  
29. treatment, vaccine, find, interesting, base, study, research, cure, scientist, science-> Vaccine  
30. pandemic, crisis, global, economy, impact, economic, future, opportunity, industry, threat->Economy threat  
31. spread, stop, child, school, close, kid, student, parent, risk, slow->Closing of Schools In UK  
32. run, walk, city, hour, open, exercise, drive, sit, close, space->Transportation/ Exercise  
33. food, order, deliver, eat, delivery, drink, pick, animal, local, restaurant->Food  
34. test, symptom, positive, show, testing, covid, result, contact, cough, flu->Test results/ Symptoms  
35. covid, high, risk, death, infection, number, rate, low, population, datum->  
36. free, update, late, information, link, check, visit, resource, advice, online-> Information links  
37. worker, fight, covid, staff, continue, protect, essential, resident, service, key-> Medical staff appreciation  
38. world, country, pandemic, leader, leadership, problem, real, show, poor, rest->Leadership opinions(Good/ Bad)

No guarantee that lower-dimensional projections capture a real-world  
relationship

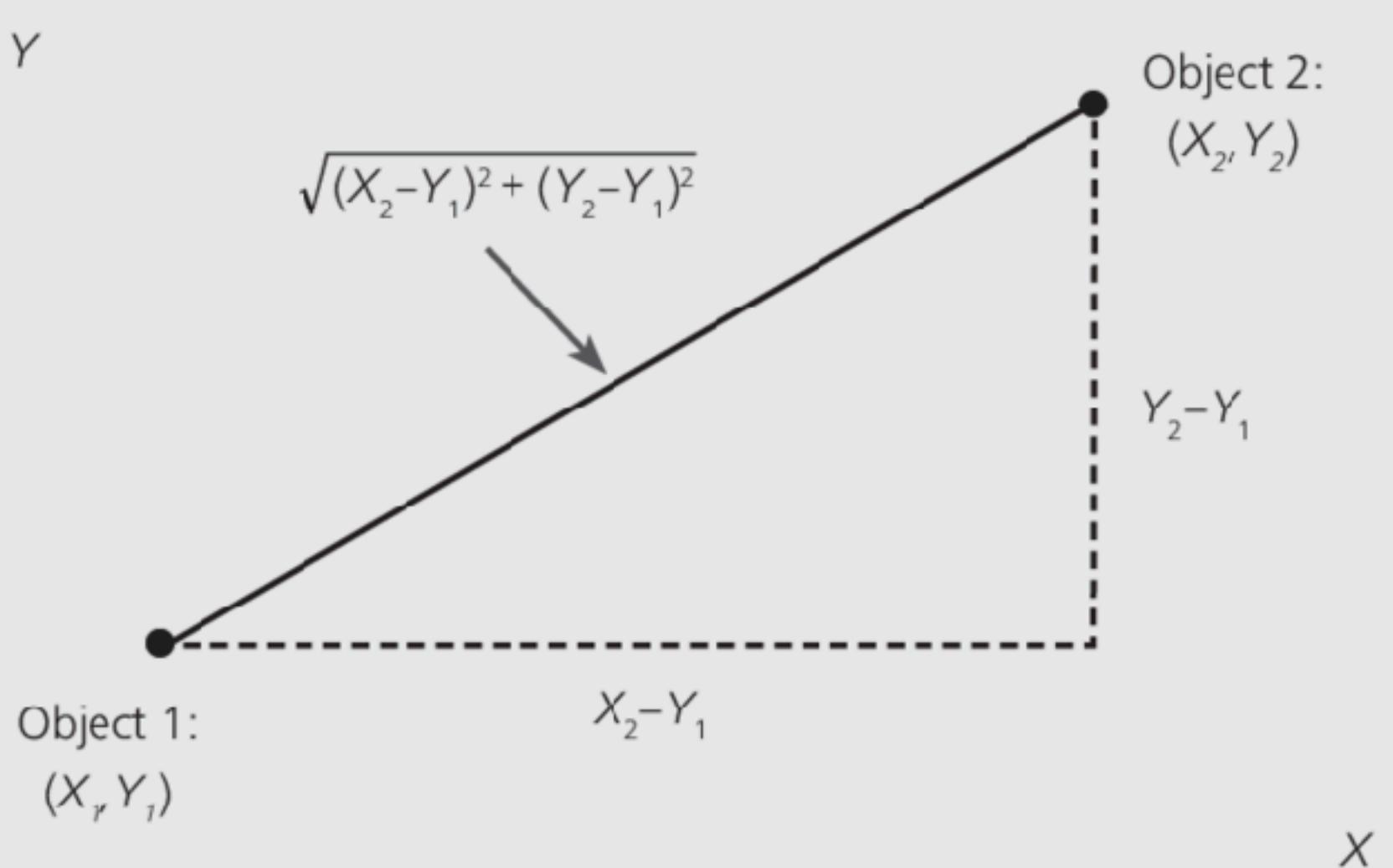
Why might we want to reduce a dataset into fewer dimensions?



Euclidean Distance  $d(x,y)$ :

Similarity:  $1/d(r_A, r_B)$

Similarity  $\rightarrow \infty$  as  $d(r_A, r_B) \rightarrow 0$



# review articles

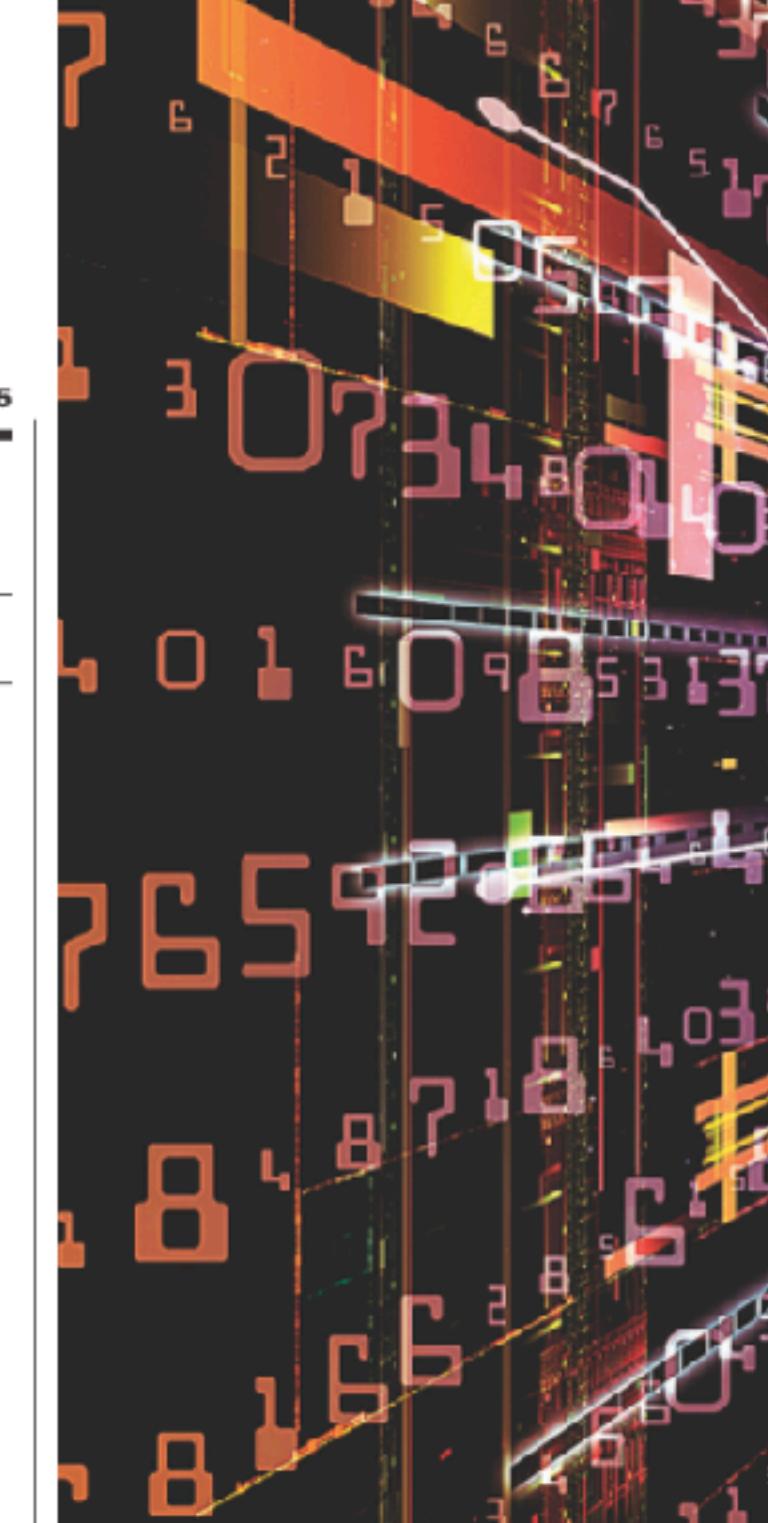
DOI:10.1145/2347736.2347755

**Tapping into the “folk knowledge” needed to advance machine learning applications.**

BY PEDRO DOMINGOS

## A Few Useful Things to Know About Machine Learning

MACHINE LEARNING SYSTEMS automatically learn programs from data. This is often a very attractive alternative to manually programming them and in the



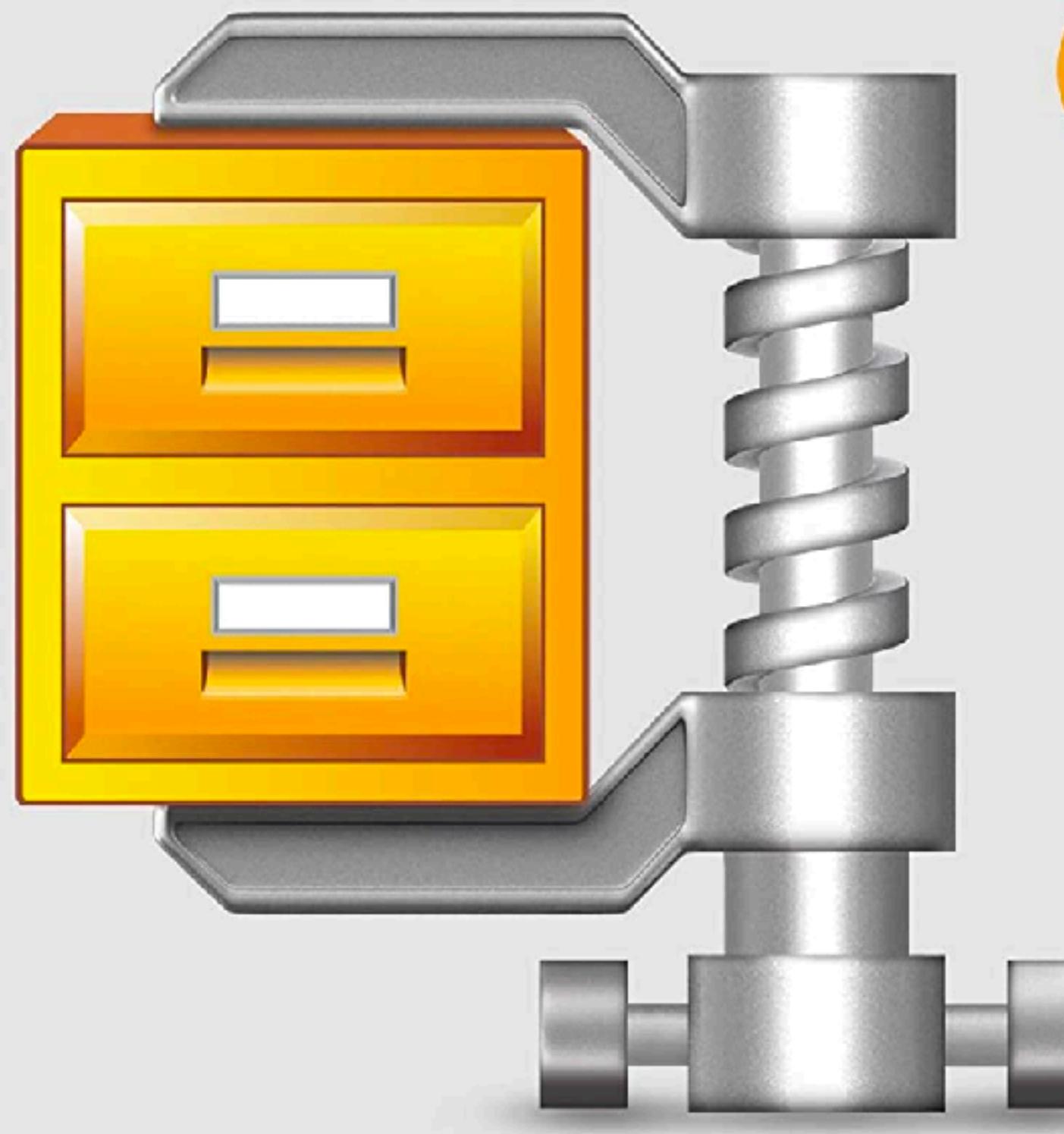
is needed to successfully develop machine learning applications is not readily available in them. As a result, many machine learning projects take much longer than necessary or wind up producing less-than-ideal results. Yet much of this folk knowledge is fairly easy to communicate. This is the purpose of this article.

» **key insights**

Uninformative for points in higher dimensions

fraud detection, stock trading, drug design, and many other applications. A recent report from the McKinsey Global Institute asserts that machine learning (a.k.a. data mining or predictive analytics) will be the driver of the next big wave of innovation.<sup>15</sup> Several fine textbooks are available to interested practitioners and researchers (for example, Mitchell<sup>16</sup> and Witten et al.<sup>24</sup>). However, much of the “folk knowledge” that

- Machine learning is widely used in computer science and other fields. However, developing successful machine learning applications requires a substantial amount of “black art” that is difficult to find in textbooks.
- This article summarizes 12 key lessons that machine learning researchers and practitioners have learned. These include pitfalls to avoid, important issues to focus on, and answers to common questions.



WinZip 8  
Mac

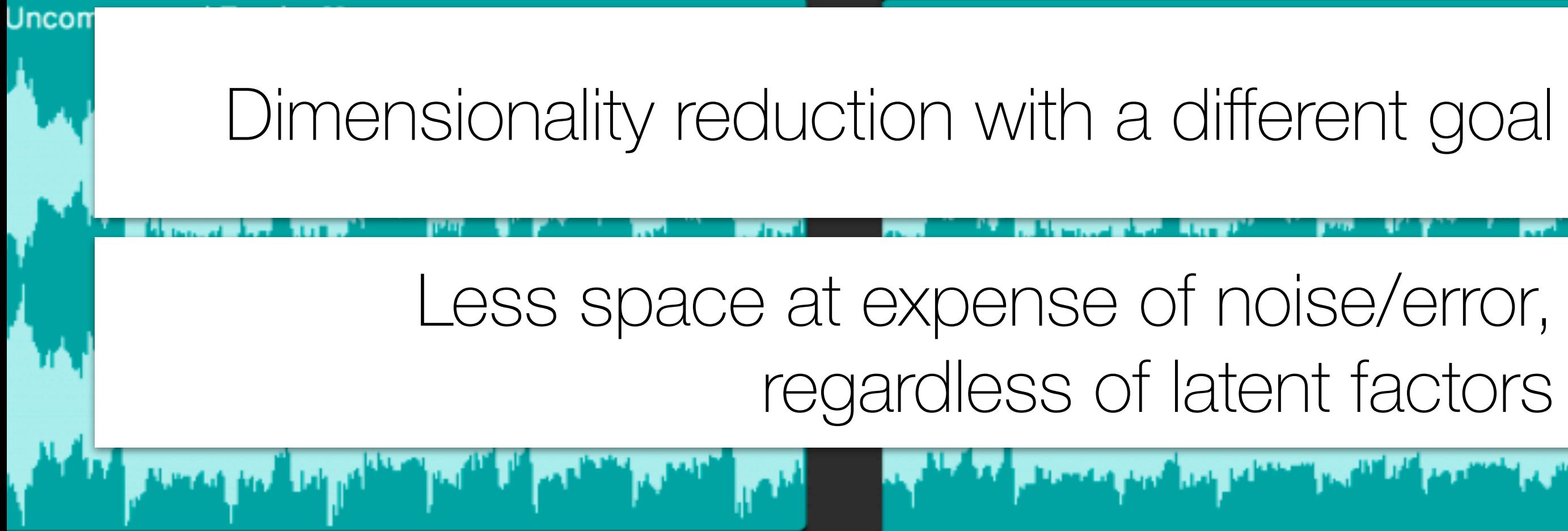


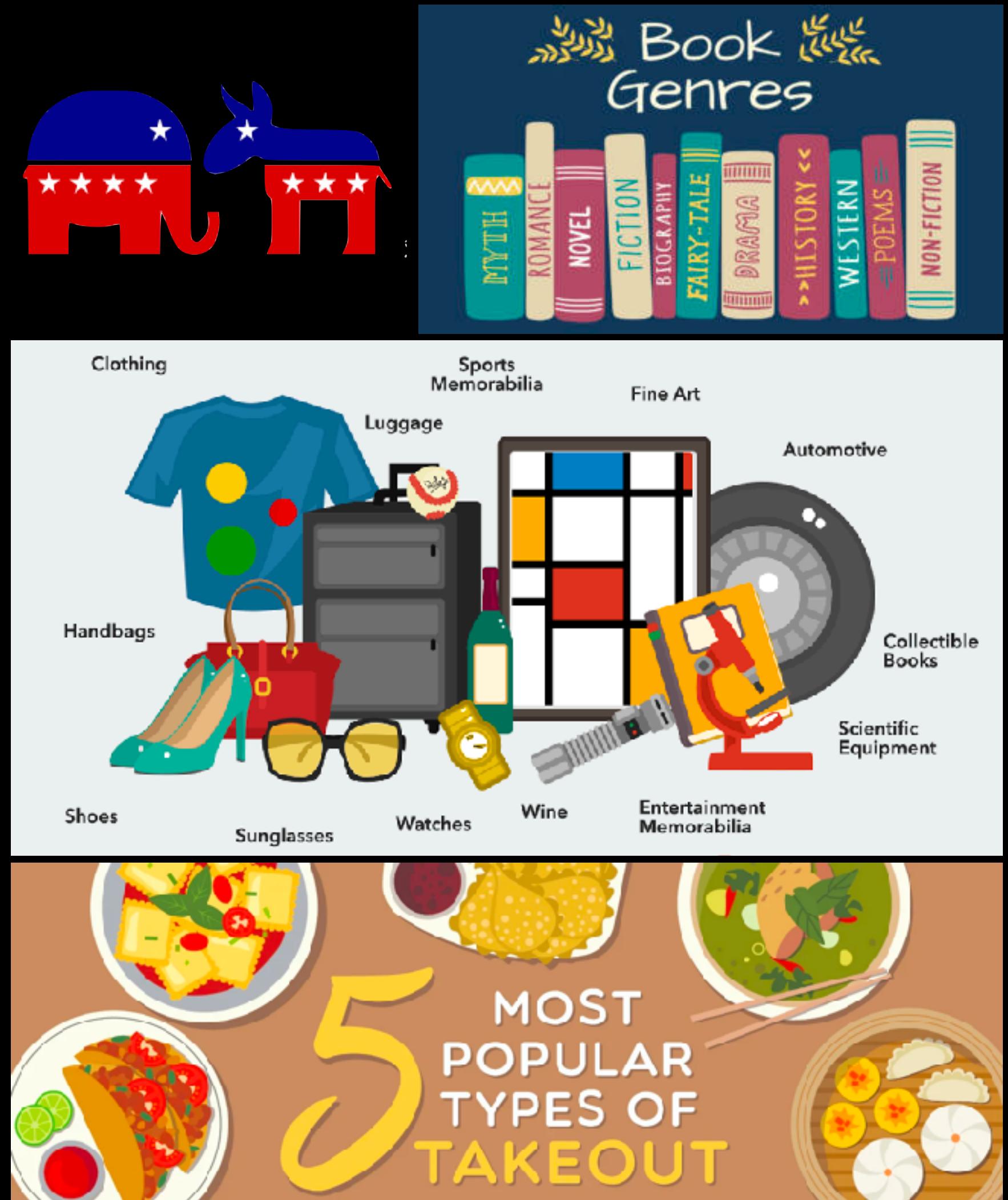
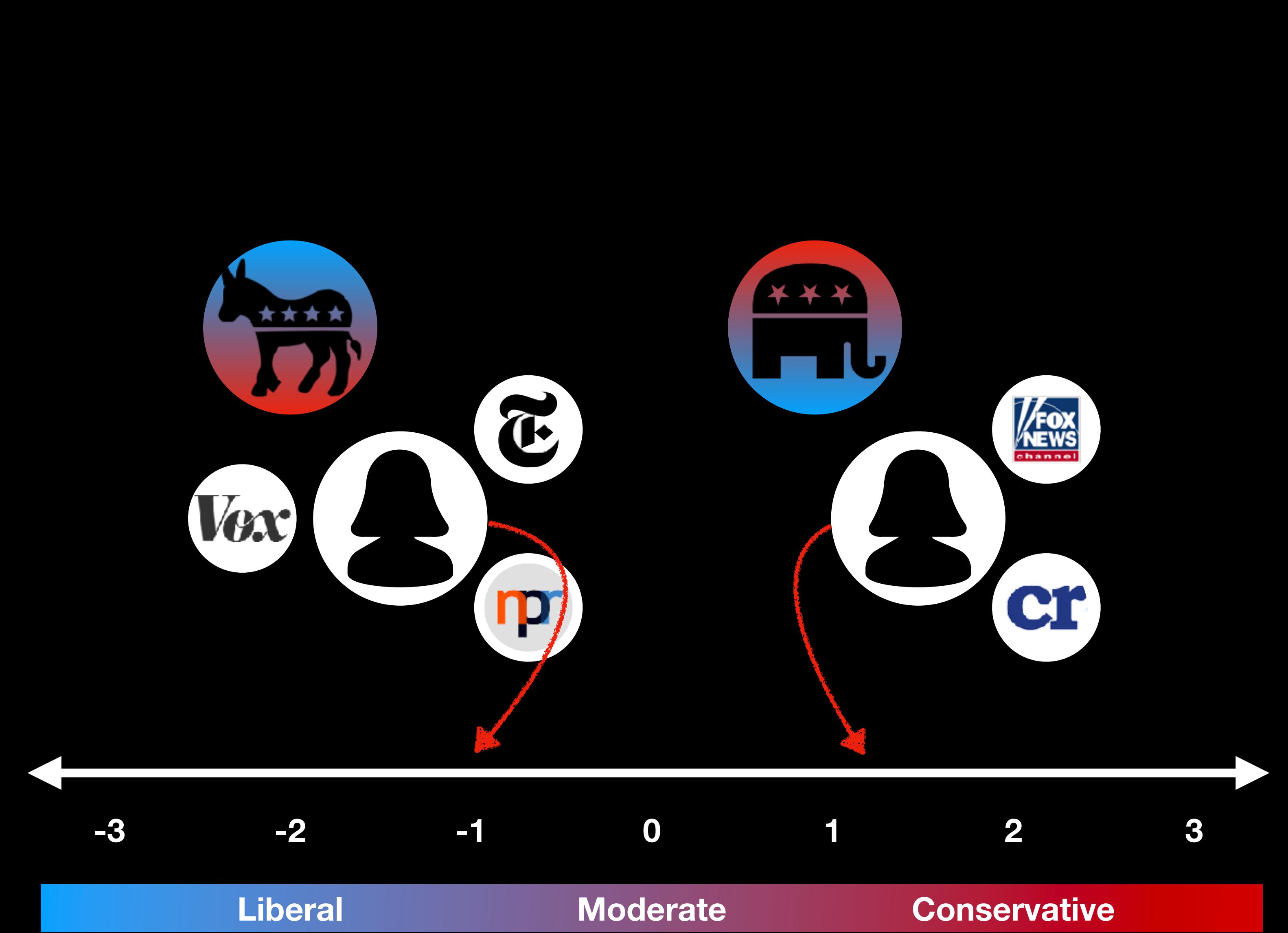
120KB  
Uncompressed → 68KB  
Compressed

## What Is Image Compression?

Uncompressed  
Waveform

Compressed  
Waveform

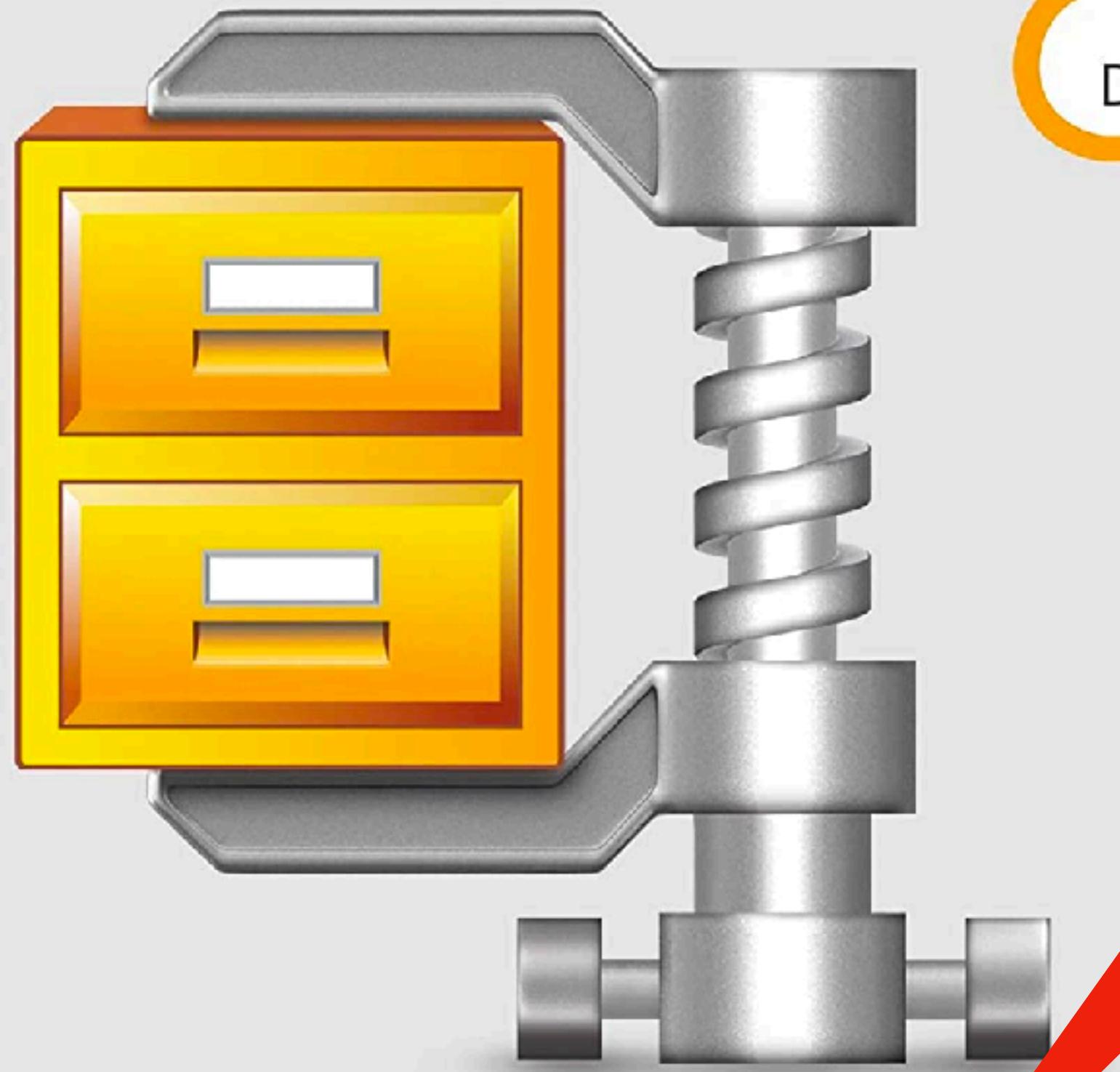




All recover some form of underlying structure

Do you need to know the latent factors to compress data?

No, you do not



WinZip 8  
Mac

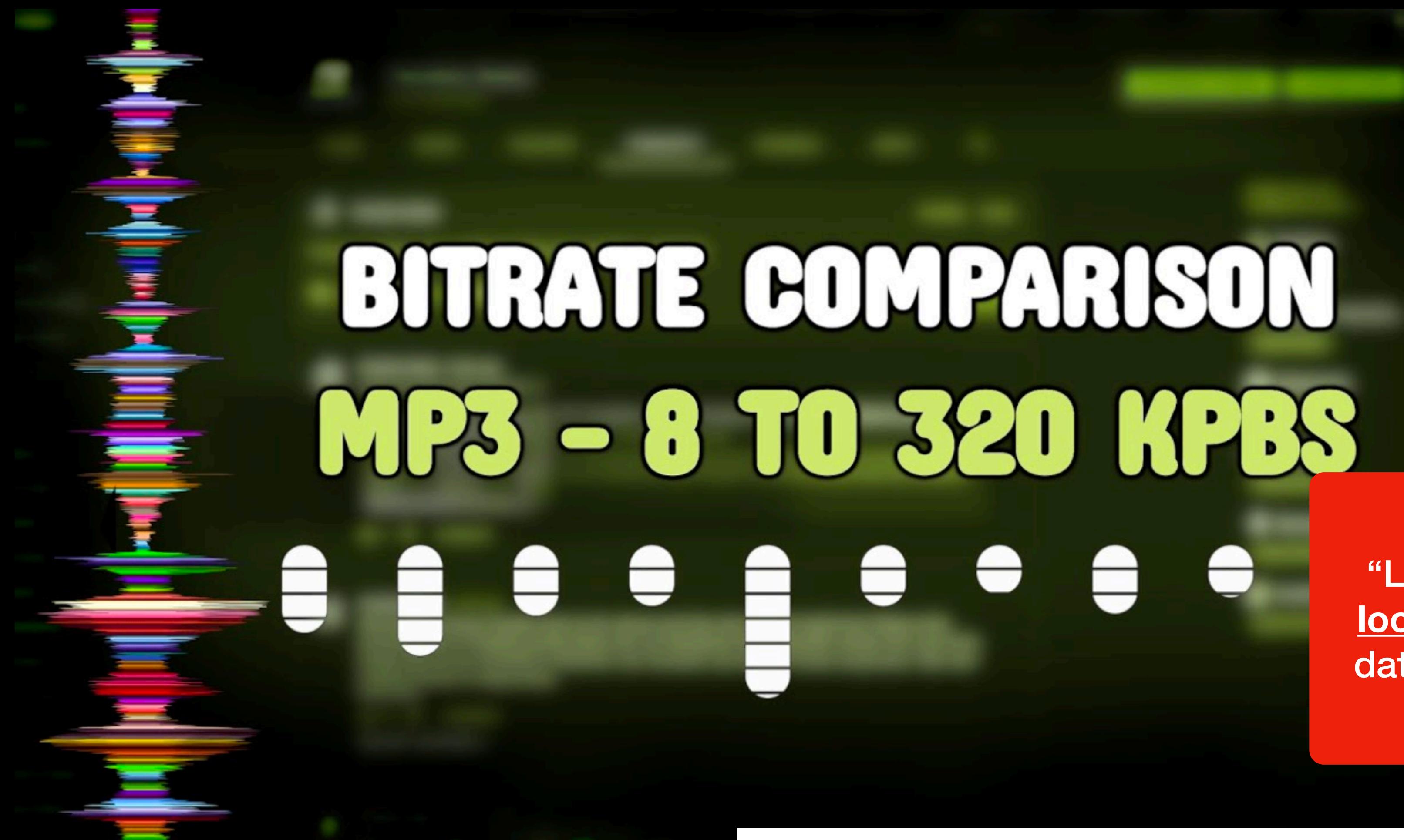


WinZip is a “lossless” compression scheme



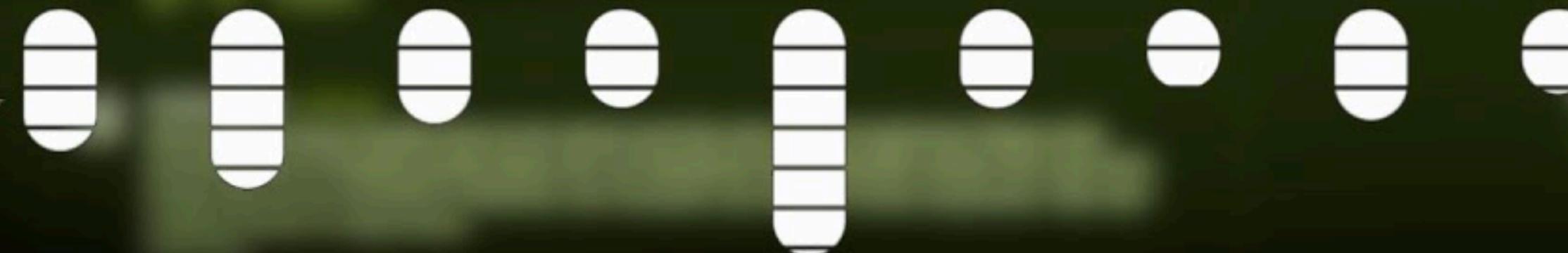
Don't need to know underlying structure

Only need to recover the “original” data



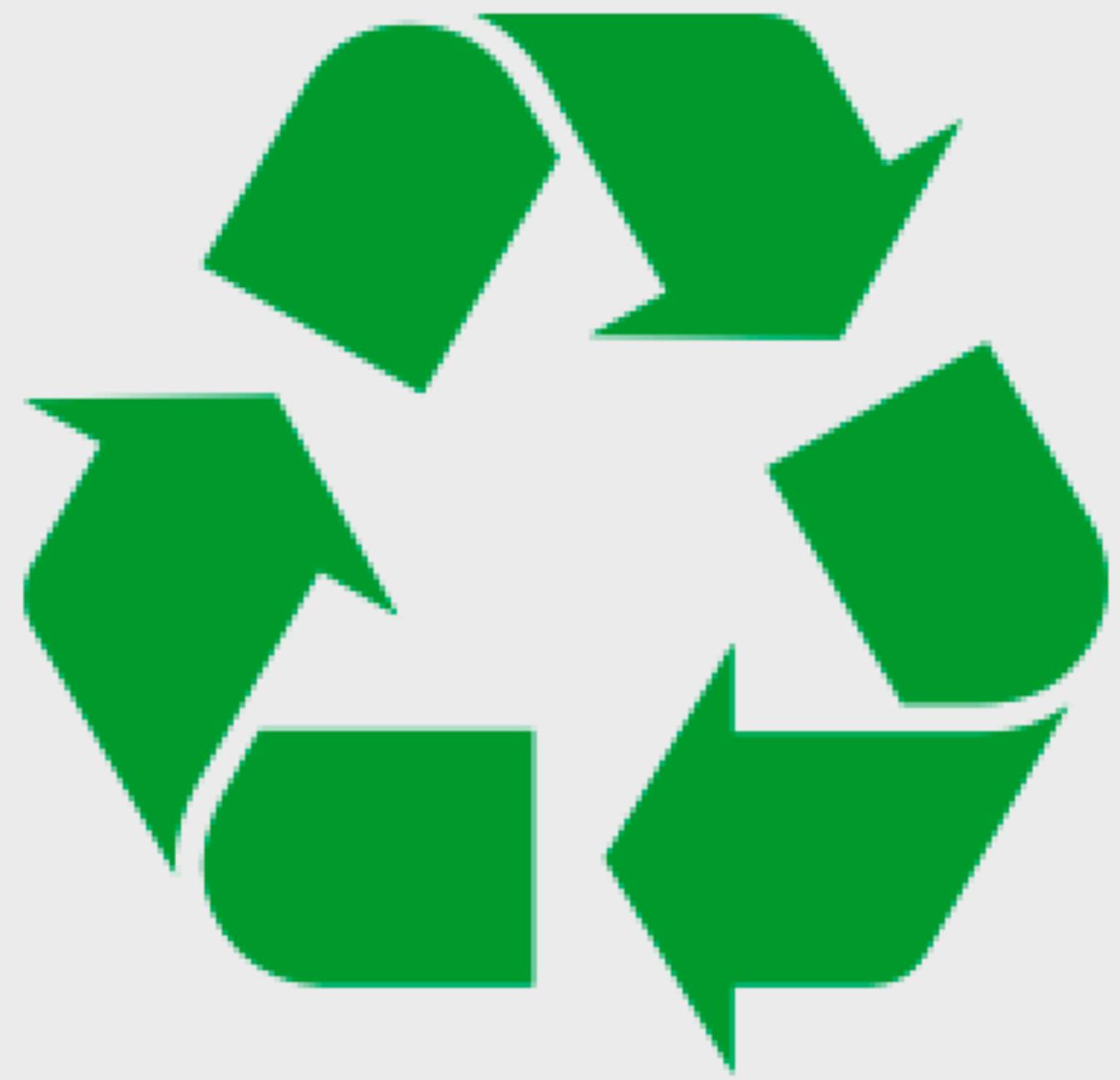
# BITRATE COMPARISON

## MP3 - 8 TO 320 KPBS



“Lossy” means we  
lose some original  
data to decompress

MP3 compresses audio in a “lossy” way



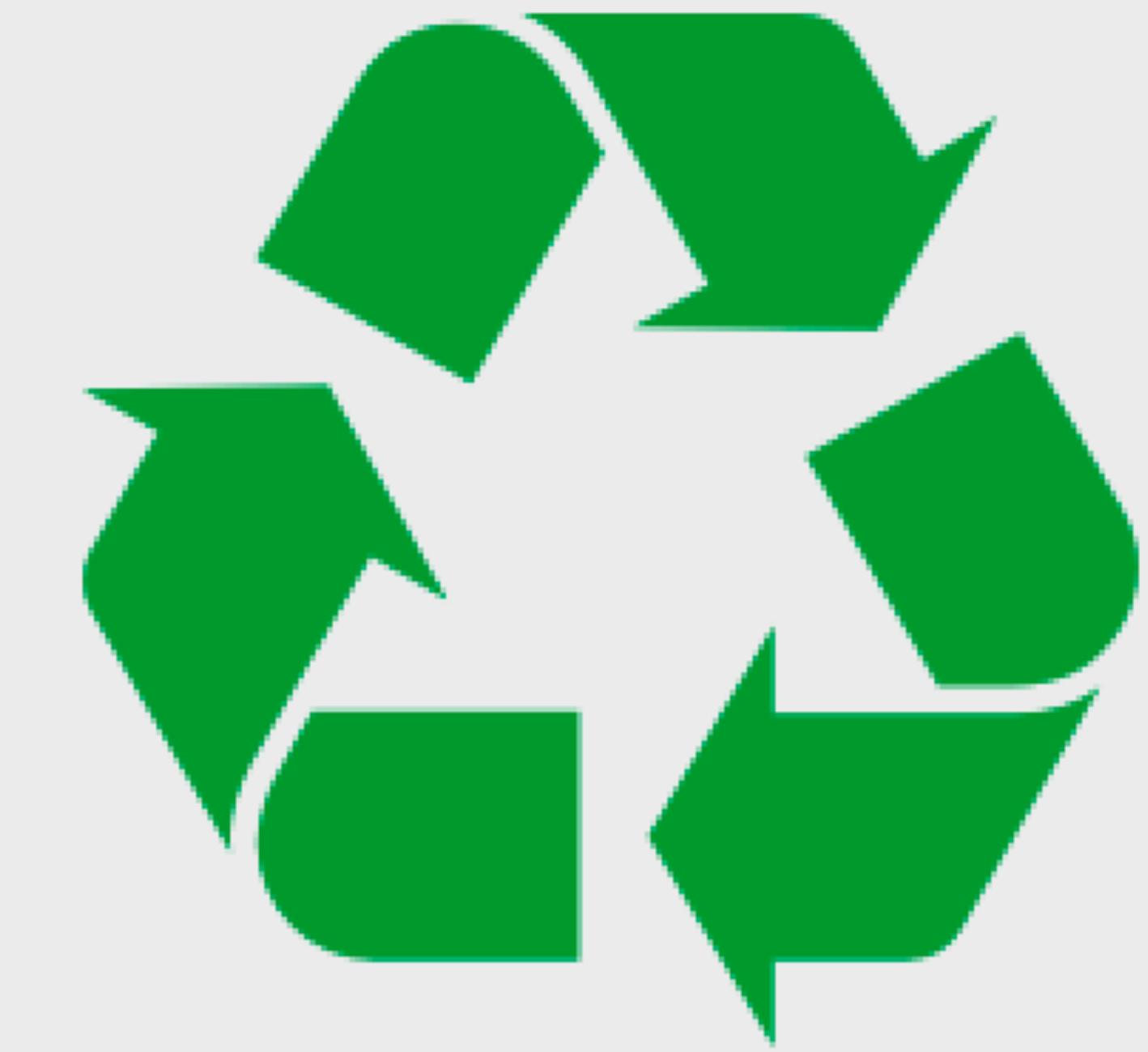
3 KB  
**PNG**



JPEG is lossy

3 KB  
**JPEG**

Quality 20



24 KB  
**JPEG**

Quality 100

Compression prioritizes reduced size while minimizing “reconstruction error”

# This Module's Learning Objectives

Define the connection between latent factors and dimensionality reduction

Describe at least two methods for dimensionality reduction

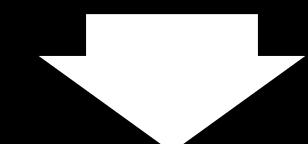
# This Module's Learning Objectives

Define the connection between latent factors and dimensionality reduction

Describe at least two methods for dimensionality reduction

**Adjacency Matrix  $A$**

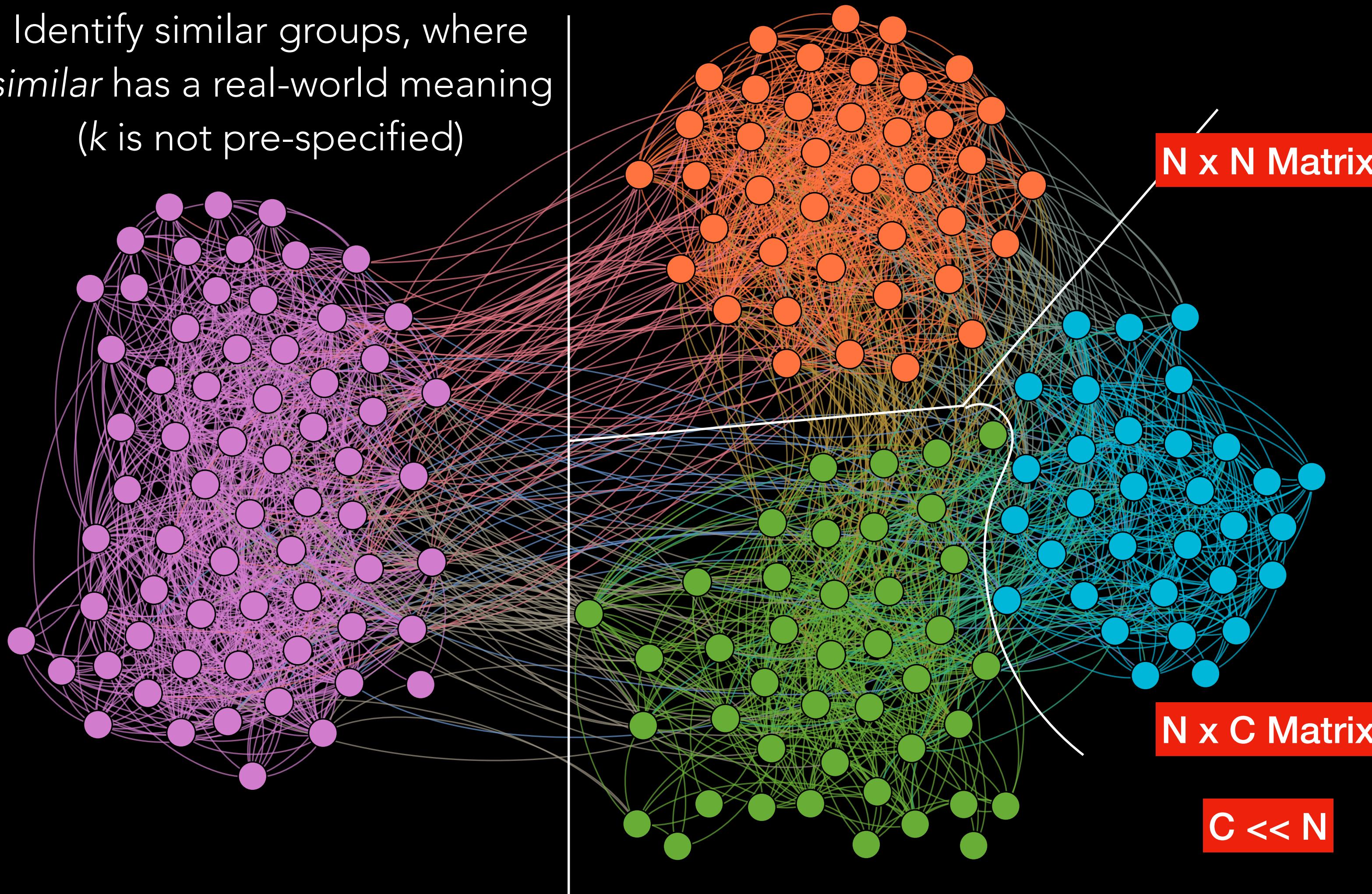
	$V_1$	$V_2$	$V_3$	$\dots$	$V_n$
$V_1$	0	0	1		1
$V_2$	1	0	0		0
$V_3$	0	1	0		0
$\dots$					
$V_n$	0	0	0		0



	$C_1$	$C_2$	$C_3$	$C_4$
$V_1$	0	0	1	0
$V_2$	1	0	0	0
$V_3$	0	1	0	0
$\dots$				
$V_n$	0	0	0	1

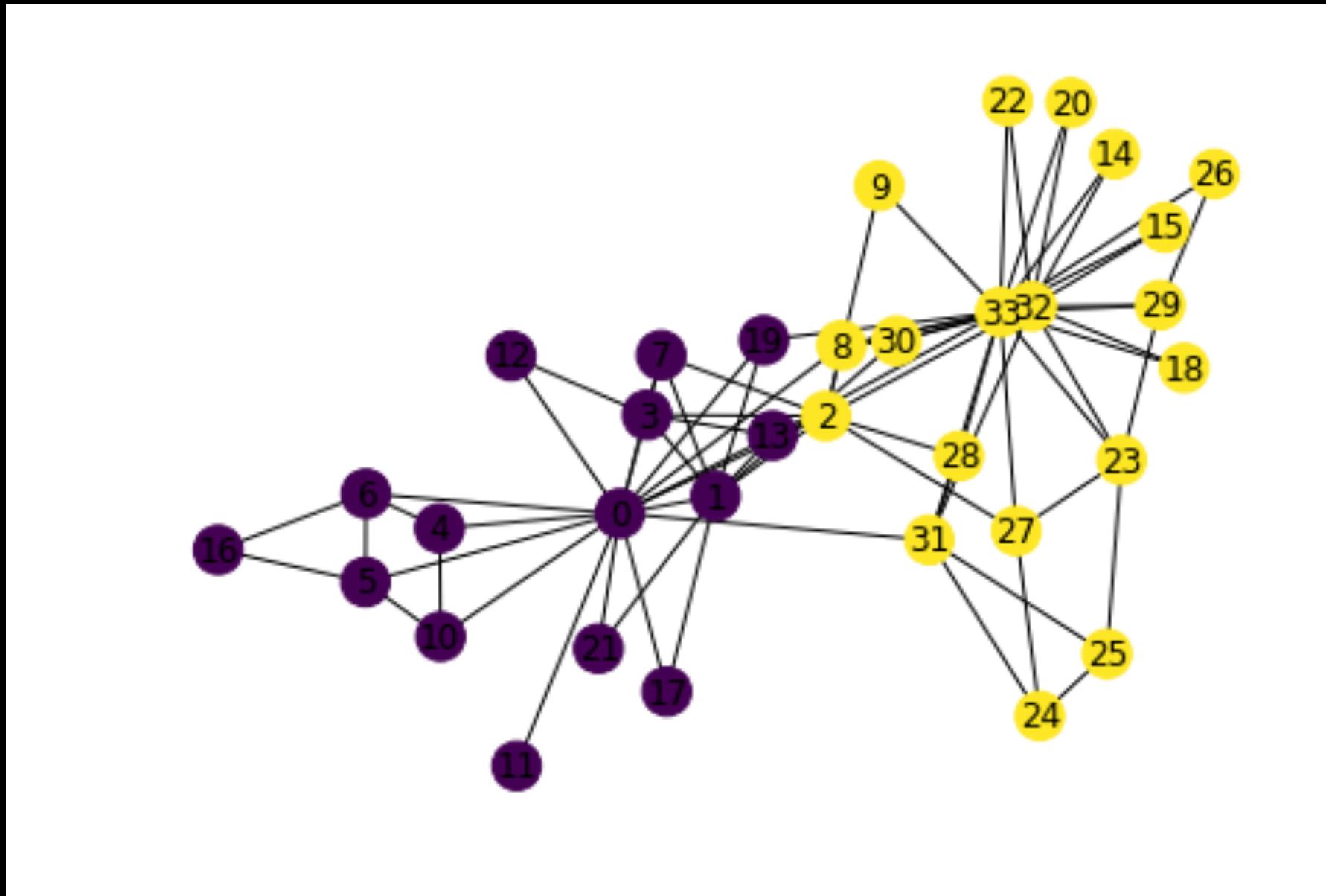
## Clustering:

Identify similar groups, where  
*similar* has a real-world meaning  
( $k$  is not pre-specified)



# Girvan-Newman Clustering

## Zachary's Karate Club Graph



Have discussed methods for extracting communities from graphs or  
selecting important features

Many general methods also exist

**Some feature matrix A**

	$d_1$	$d_2$	$d_3$	$\dots$	$d_n$
$x_1$	0	0	1		1
$x_2$	1	0	0		0
$x_3$	0	1	0		0
$\dots$					
$x_m$	0	0	0		0

Decompose to latent factors

Reduce dimensions with minimum  
reconstruction error

Some feature matrix A

	$d_1$	$d_2$	$d_3$	$\dots$	$d_n$
$x_1$	0	0	1		1
$x_2$	1	0	0		0
$x_3$	0	1	0		0
$\dots$					
$x_m$	0	0	0		0



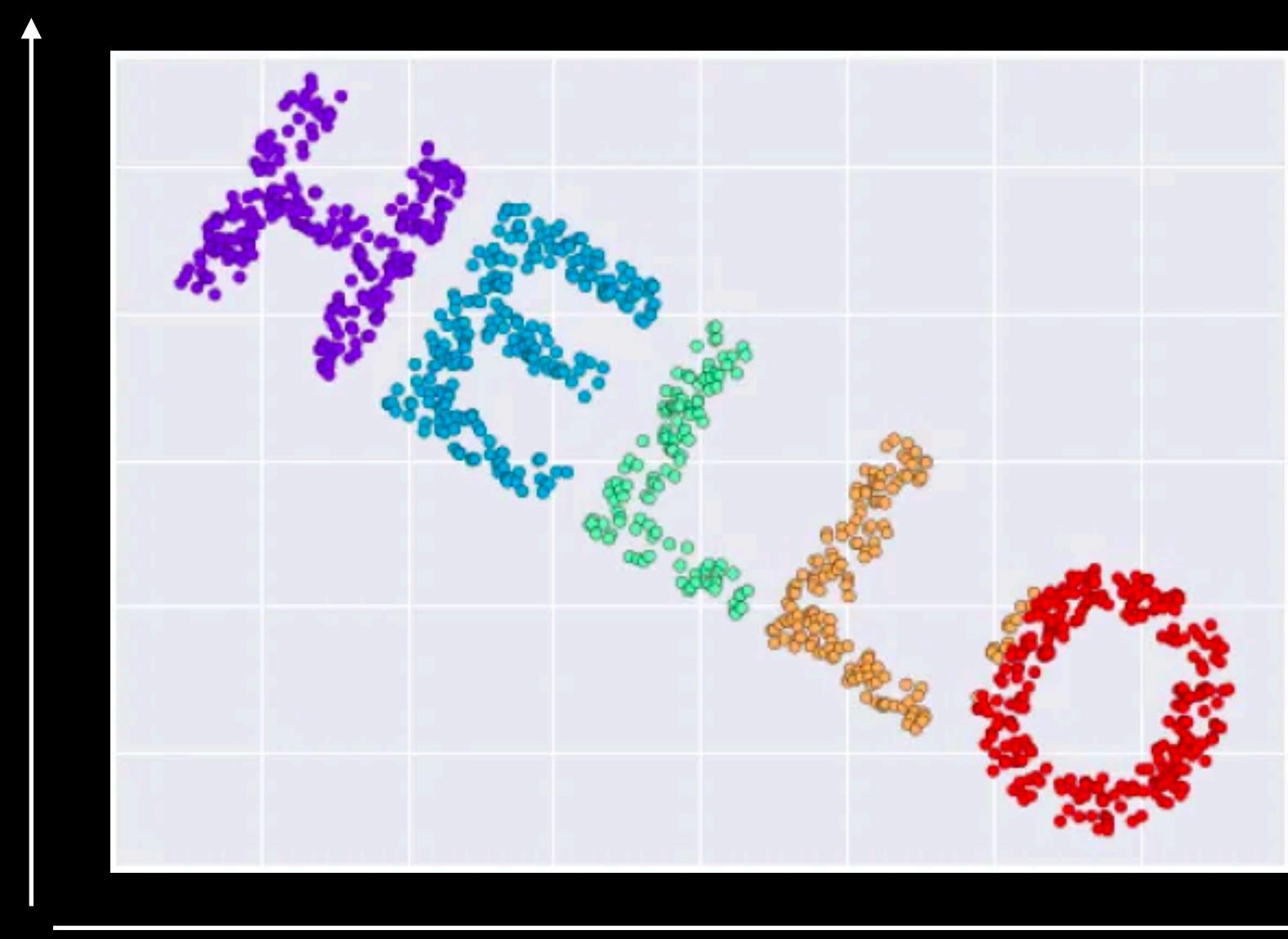
First Principal Component  
Second Principal Component

	$d'_1$	$d'_2$
$x_1$	0	0
$x_2$	1	0
$x_3$	0	1
$\dots$		
$x_n$	0	0

n components ordered by eigenvalues

Keeping all components produces a noiseless rotation

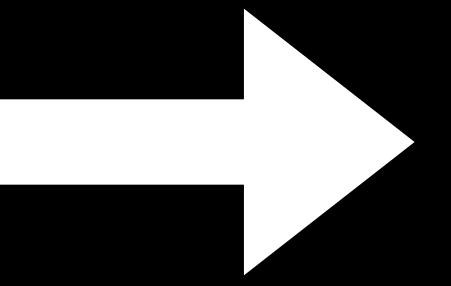
Principal Component Analysis (PCA)



Project onto dimension with highest variance

Can achieve the same reduction through a process of “decomposition”

**Some feature matrix  $A$  ( $m \times n$ )**



**Sub-matrix  $U$  ( $m \times d$ )**

$\times$

**Sub-matrix  $V$  ( $d \times n$ )**

Or matrix factorization

	Feature 1	Feature 2
User 1	0.24	0.16
User 2	0.28	0.18
User 3	0.3	0.2
User 4	0.24	0.16

**User Matrix**

	Feature 1	Feature 2
User 1	0.24	0.16
User 2	0.28	0.18
User 3	0.3	0.2
User 4	0.24	0.16

**User Matrix**

X

	1.5	1.2	1.0	0.8
	1.7	0.6	1.1	0.4

**Home Matrix**

	Feature 1	Feature 2
User 1	0.24	0.16
User 2	0.28	0.18
User 3	0.3	0.2
User 4	0.24	0.16

**User Matrix**

X

	House 1	House 2	House 3	House 4
Feature 1	1.5	1.2	1.0	0.8
Feature 2	1.7	0.6	1.1	0.4

**Home Matrix**

	House 1	House 2	House 3	House 4
User 1	0.2	0.9	0.4	0.2
User 2	0.8	0.2	0.7	0.2
User 3	0.2	1.0	0.2	0.4
User 4	0.2	0.7	0.8	0.2

But what if you could go backwards?

				
	0.2	<b>0.9</b>	<b>0.4</b>	0.2
	<b>0.8</b>	0.2	<b>0.7</b>	0.2
	0.2	<b>1.0</b>	0.2	<b>0.4</b>
	0.2	<b>0.7</b>	<b>0.8</b>	0.2

**User-Home Matrix**

				
	0.2	<b>0.9</b>	<b>0.4</b>	0.2
	<b>0.8</b>	0.2	<b>0.7</b>	0.2
	0.2	<b>1.0</b>	0.2	<b>0.4</b>
	0.2	<b>0.7</b>	<b>0.8</b>	0.2

**User-Home Matrix**

≡

What approximately equals this matrix?

	0.2	0.9	0.4	0.2
	0.8	0.2	0.7	0.2
	0.2	1.0	0.2	0.4
	0.2	0.7	0.8	0.2

User-Home Matrix

≡

	Feature 1	Feature 2
	0.24	0.16
	0.28	0.18
	0.3	0.2
	0.24	0.16

**User Matrix**

	0.2	0.9	0.4	0.2
	0.8	0.2	0.7	0.2
	0.2	1.0	0.2	0.4
	0.2	0.7	0.8	0.2

User-Home Matrix

	Feature 1	Feature 2
	0.24	0.16
	0.28	0.18
	0.3	0.2
	0.24	0.16

**User Matrix**

Feature 1	1.5	1.2	1.0	0.8
Feature 2	1.7	0.6	1.1	0.4
	<b>Home Matrix</b>			

	0.2	<b>0.9</b>	0.4	0.0
	<b>0.8</b>	0.2		
	0.2	1.0	0.2	0.4
	0.2	0.7	0.8	0.2

User-Home Matrix

Reduced feature set

	Feature 1	Feature 2
	<b>0.24</b>	0.16
	0.28	0.18
	0.3	0.2
	0.24	0.16

User Matrix

Feature 1	1.5	1.2	1.0	0.8
Feature 2	1.7	0.6	1.1	0.4

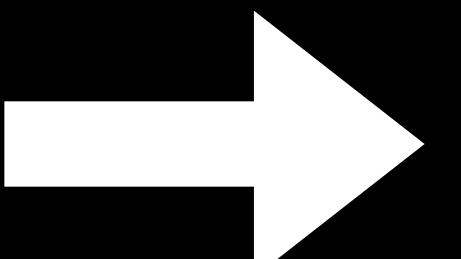
Home Matrix

Some feature matrix  $A$  ( $m \times n$ )

	$d_1$	$d_2$	$d_3$	$\dots$	$d_n$
$x_1$	0	0	1		1
$x_2$	1	0	0		0
$x_3$	0	1	0		0
$\dots$					
$x_m$	0	0	0		0

Matrix  $U$  ( $m \times r$ )

	$d'_1$	$\dots$	$d'_r$
$x_1$	0		0
$x_2$	1		0
$x_3$	0		1
$\dots$			
$x_m$	0		0



$\times$

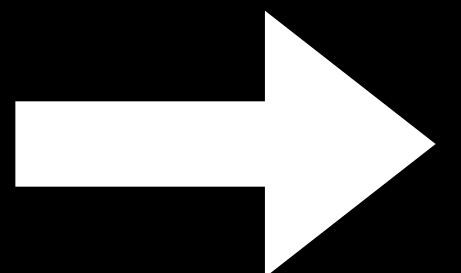
Matrix  $V$  ( $n \times r$ )

	$d'_1$	$\dots$	$d'_r$
$d_1$	0		0
$d_2$	1		0
$d_3$	0		1
$\dots$			
$d_n$	0		0

$T$

Some feature matrix  $A$  ( $m \times n$ )

	$d_1$	$d_2$	$d_3$	$\dots$	$d_n$
$x_1$	0	0	1		1
$x_2$	1	0	0		0
$x_3$	0	1	0		0
$\dots$					
$x_m$	0	0	0		0



Matrix  $U$  ( $m \times r$ )

	$d'_1$	$\dots$	$d'_r$
$x_1$	0		0
$x_2$	1		0
$x_3$	0		1
$\dots$			
$x_m$	0		0

Matrix  $\Sigma$  ( $r \times r$ )

$s_1$	0	0		0
0	$s_2$	0		0
0	0	$s_3$		0
			$\dots$	
0	0	0	0	$s_r$

Matrix  $V$  ( $n \times r$ )

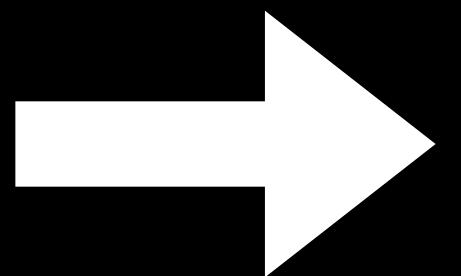
	$d'_1$	$\dots$	$d'_r$
$d_1$	0		0
$d_2$	1		0
$d_3$	0		1
$\dots$			
$d_n$	0		0

$T$

Singular Value Decomposition (SVD)

Some feature matrix  $A$

	$d_1$	$d_2$	$d_3$	$\dots$	$d_n$
$x_1$	0	0	1		1
$x_2$	1	0	0		0
$x_3$	0	1	0		0
$\dots$					
$x_m$	0	0	0		0



Matrix  $U$  ( $n \times r$ )

	$d'_1$	$\dots$	$d'_r$
$x_1$	0		0
$x_2$	1		0
$x_3$	0		1
$\dots$			
$x_n$	0		0

Matrix  $\Sigma$  ( $r \times r$ )

$s_1$	0	0	0
0	$s_2$	0	0
0	0	$s_3$	0
0	0	0	0

$\times$

Matrix  $V$  ( $n \times r$ )

	$d'_1$	$\dots$	$d'_r$
$v_1$	0		0
$v_2$	1		0
$v_3$	0		1
$\dots$			
$v_n$	0		0

$T$

Keep only important components

# SVD - Definition

$$A_{[m \times n]} = U_{[m \times r]} \Sigma_{[r \times r]} (V_{[n \times r]})^T$$

## ■ A: Input data matrix

- $m \times n$  matrix (e.g.,  $m$  documents,  $n$  terms)

## ■ U: Left singular vectors

- $m \times r$  matrix ( $m$  documents,  $r$  concepts)

## ■ $\Sigma$ : Singular values

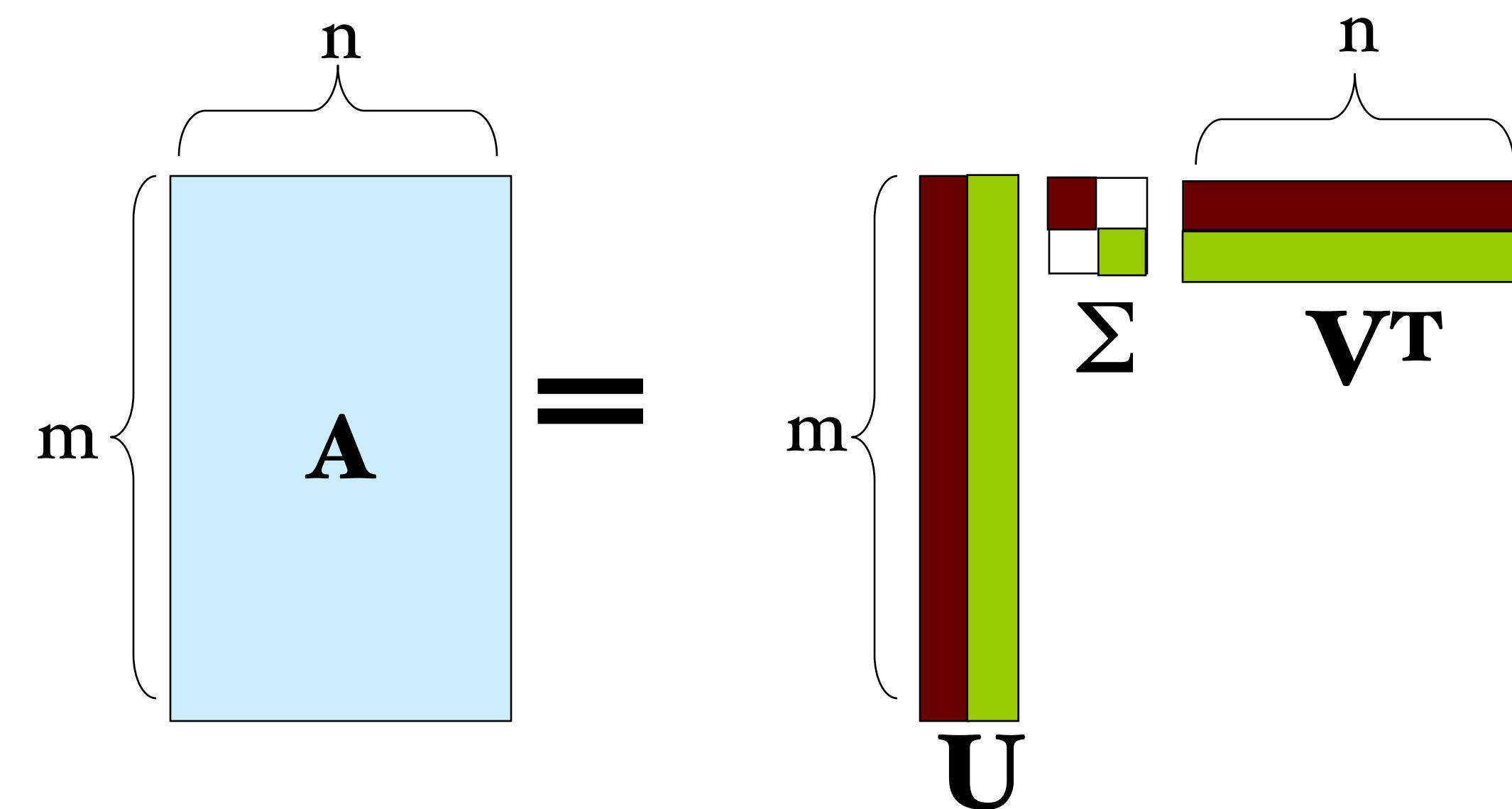
- $r \times r$  diagonal matrix (strength of each ‘concept’)

## ■ V: Right singular vectors

- $n \times r$  matrix ( $n$  terms,  $r$  concepts)

# SVD

$$\mathbf{A} \approx \mathbf{U}\Sigma\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i^\top$$



# SVD - Properties

It is **always** possible to decompose a real matrix  $A$  into  $A = U \Sigma V^T$ , where

- $U, \Sigma, V$ : unique

- $\Sigma$ : diagonal

- Entries (**singular values**) are **positive**,  
and sorted in decreasing order ( $s_1 \geq s_2 \geq \dots \geq 0$ )

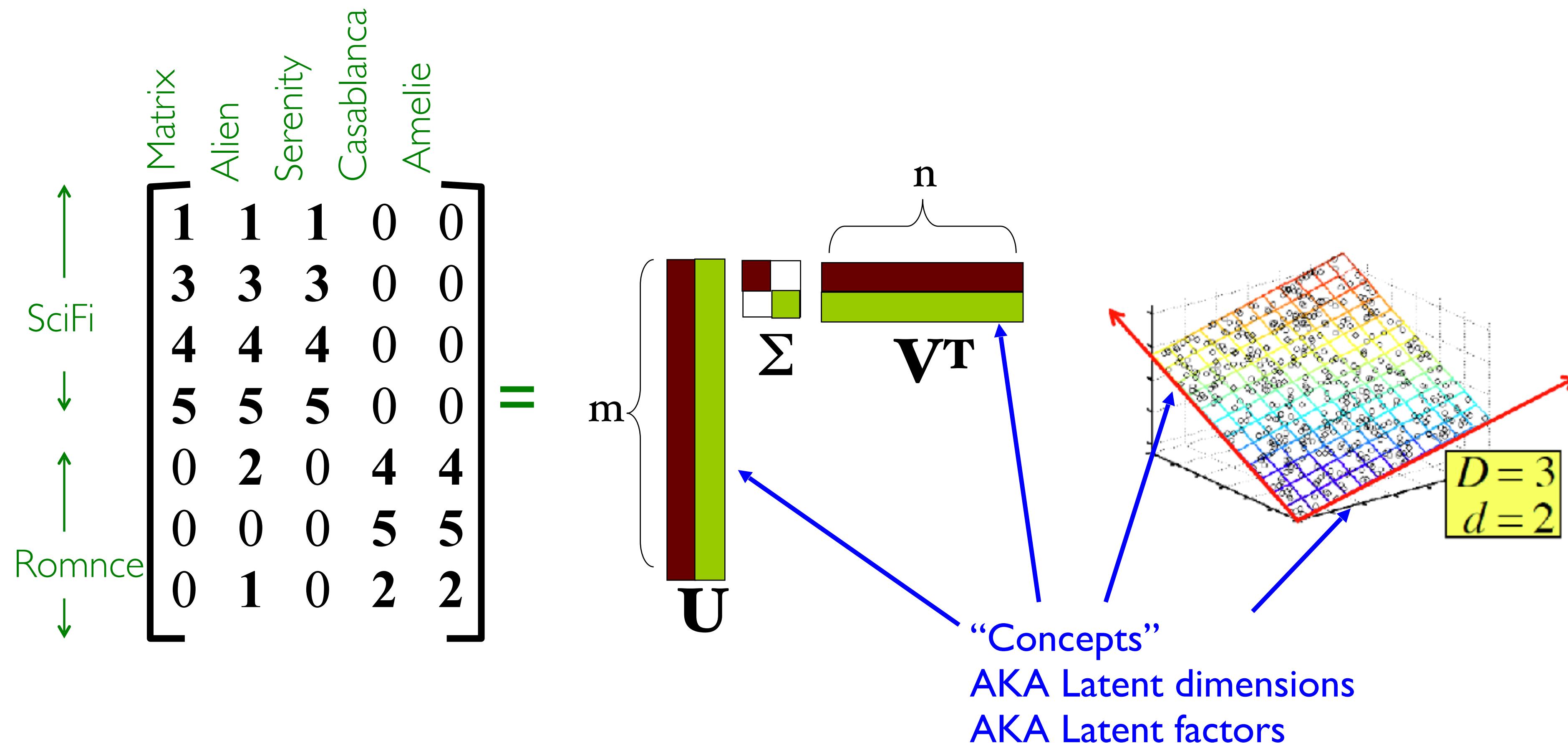
$$\Sigma =$$

s1	0	0		0
0	s2	0		0
0	0	s3		0
			...	
0	0	0	0	sr

Nice proof of uniqueness: <http://www.mpi-inf.mpg.de/~bast/ir-seminar-ws04/lecture2.pdf>

# SVD – Example: Users-to-Movies

■  $A = U \Sigma V^T$  - example: Users to Movies



# SVD – Example: Users-to-Movies

■  $A = U \Sigma V^T$  - example: Users to Movies

$$\begin{matrix} & \text{Matrix} \\ \uparrow & \begin{matrix} \text{Alien} & \text{Serenity} & \text{Casablanca} & \text{Amelie} \end{matrix} \\ \text{SciFi} & \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \\ \downarrow & \\ \uparrow & \begin{matrix} \text{Romance} \end{matrix} \\ \downarrow & \end{matrix} = \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

# SVD – Example: Users-to-Movies

■  $A = U \Sigma V^T$  - example: Users to Movies

$$\begin{array}{c}
 \text{Matrix} \\
 \begin{array}{c}
 \begin{array}{c} \uparrow \\ \text{SciFi} \\ \downarrow \\ \uparrow \\ \text{Romnce} \\ \downarrow \end{array}
 \left[ \begin{array}{ccccc}
 1 & 1 & 1 & 0 & 0 \\
 3 & 3 & 3 & 0 & 0 \\
 4 & 4 & 4 & 0 & 0 \\
 5 & 5 & 5 & 0 & 0 \\
 0 & 2 & 0 & 4 & 4 \\
 0 & 0 & 0 & 5 & 5 \\
 0 & 1 & 0 & 2 & 2
 \end{array} \right]
 \end{array}
 \end{array}
 = \left[ \begin{array}{ccc}
 0.13 & 0.02 & -0.01 \\
 0.41 & 0.07 & -0.03 \\
 0.55 & 0.09 & -0.04 \\
 0.68 & 0.11 & -0.05 \\
 0.15 & -0.59 & 0.65 \\
 0.07 & -0.73 & -0.67 \\
 0.07 & -0.29 & 0.32
 \end{array} \right]
 \times \left[ \begin{array}{ccc}
 12.4 & 0 & 0 \\
 0 & 9.5 & 0 \\
 0 & 0 & 1.3
 \end{array} \right]
 \times \left[ \begin{array}{ccccc}
 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\
 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\
 0.40 & -0.80 & 0.40 & 0.09 & 0.09
 \end{array} \right]$$

SciFi-concept      Romance-concept

# SVD – Example: Users-to-Movies

■  $A = U \Sigma V^T$  - example:

$U$  is “user-to-concept”  
similarity matrix

Matrix

	Alien	Serenity	Casablanca	Amelie
SciFi	1 1 1 0 0	0.13 0.02 -0.01	0.41 0.07 -0.03	0.55 0.09 -0.04
	3 3 3 0 0	0.68 0.11 -0.05	0.15 -0.59 0.65	0.07 -0.73 -0.67
	4 4 4 0 0	0.15 -0.59 0.65	0.07 -0.73 -0.67	0.07 -0.29 0.32
	5 5 5 0 0	0.07 -0.73 -0.67	0.40 -0.80 0.40	0.12 -0.02 0.12
Romnce	0 2 0 4 4	0.07 -0.29 0.32	0.56 0.59 0.56	0.09 0.09 0.09
	0 0 0 5 5	0.56 0.59 0.56	0.12 -0.02 0.12	0.09 0.09 0.09
	0 1 0 2 2	0.40 -0.80 0.40	-0.69 -0.69 -0.69	0.09 0.09 0.09

=

SciFi-concept	Romance-concept
0.13	0.02
0.41	0.07
0.55	0.09
0.68	0.11
0.15	-0.59
0.07	-0.73
0.07	-0.29

$\times$

12.4 0 0
0 9.5 0
0 0 1.3

$\times$

0.56 0.59 0.56 0.09 0.09
0.12 -0.02 0.12 -0.69 -0.69
0.40 -0.80 0.40 0.09 0.09

# SVD – Example: Users-to-Movies

## ■ $A = U \Sigma V^T$ - example:

The diagram illustrates the calculation of the "strength" of the SciFi-concept from a movie matrix. The matrix has rows labeled by genres (SciFi, Romance) and columns labeled by movies (Matrix, Alien, Serenity, Casablanca, Amelie). The matrix is multiplied by a vector representing the SciFi-concept, resulting in a vector where the first element is circled in blue and labeled "12.4".

**Matrix**

	Matrix	Alien	Serenity	Casablanca	Amelie
SciFi	1	1	1	0	0
Romance	0	2	0	4	4
	0	0	0	5	5
	0	1	0	2	2

**SciFi-concept**

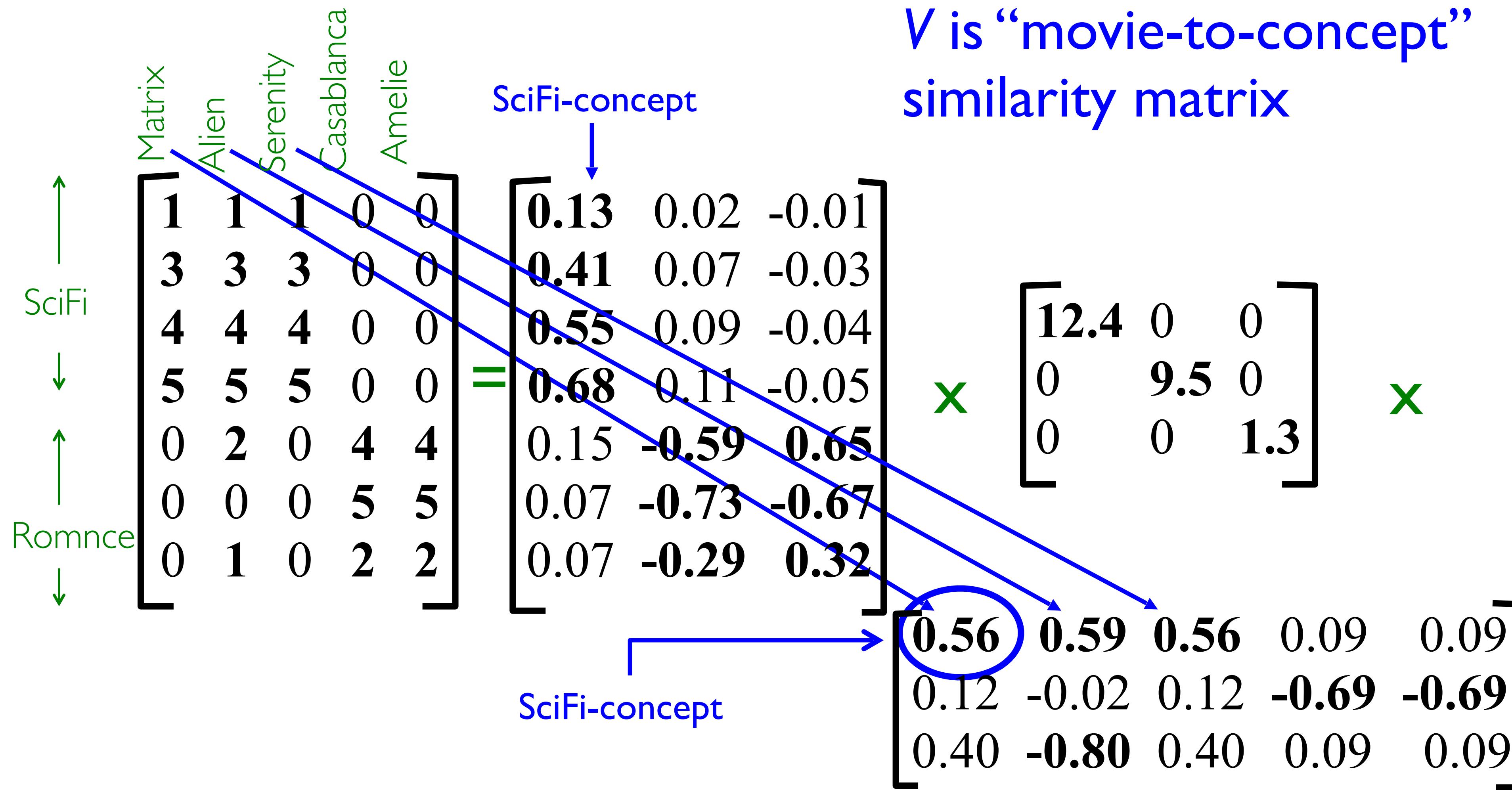
**“strength” of the SciFi-concept**

**Strength Vector:**

12.4	0	9.5	0
0	0	1.3	

# SVD – Example: Users-to-Movies

■  $A = U \Sigma V^T$  - example:



# SVD - One Interpretation

$$A_{[m \times n]} = U_{[m \times r]} \Sigma_{[r \times r]} (V_{[n \times r]})^T$$

User  
Interests      Concept  
Strengths  
(axes of most  
variation)      Movie  
Genre  
assignments

# SVD - Conclusions

## ■ SVD: $A = U \Sigma V^T$ : unique

- $U$ : user-to-concept similarities
- $V$ : movie-to-concept similarities
- $\Sigma$  : strength of each concept

## ■ Dimensionality reduction:

- keep the few largest singular values  
(80-90% of ‘energy’)
- SVD: picks up linear correlations

Prev Up Next

scikit-learn 0.24.1

[Other versions](#)

Please [cite us](#) if you use the software.

[sklearn.decomposition.TruncatedSVD](#)

Examples using

[sklearn.decomposition.TruncatedSVD](#)

## sklearn.decomposition.TruncatedSVD

```
class sklearn.decomposition.TruncatedSVD(n_components=2, *, algorithm='randomized', n_iter=5,  
random_state=None, tol=0.0)
```

[\[source\]](#)

Dimensions of output data (n\_samples, n\_components).

This transformer performs dimensionality reduction by means of truncated singular value decomposition (SVD).

Contrary to [PCA](#), it does not require that the input data be centered before computing the singular value decomposition. This means it can work on sparse matrices.

In particular, truncated SVD works on term count/tf-idf matrices as returned by the vectorizers in [sklearn.feature\\_extraction.text](#). In that context, it is known as latent semantic analysis (LSA).

This estimator supports two algorithms: a fast randomized SVD solver, and a “naive” algorithm that uses ARPACK as an eigensolver on  $X * X.T$  or  $X.T * X$ , whichever is more efficient.

Read more in the [User Guide](#).

**Parameters:**

**n\_components : int, default=2**

Desired dimensionality of output data. Must be strictly less than the number of features. The default value is useful for visualisation. For LSA, a value of 100 is recommended.

**algorithm : {'arpack', 'randomized'}, default='randomized'**

SVD solver to use. Either “arpack” for the ARPACK wrapper in SciPy (`scipy.sparse.linalg.svds`), or “randomized” for the randomized algorithm due to Halko (2009).

**n\_iter : int, default=5**

Number of iterations for randomized SVD solver. Not used by ARPACK. The default is larger than the default in [randomized\\_svd](#) to handle sparse matrices that may have large slowly decaying spectrum.

**random\_state : int, RandomState instance or None, default=None**

Used during randomized svd. Pass an int for reproducible results across multiple function calls. See

n\_components is the reduced dimensionality

# This Module's Learning Objectives

Define the connection between latent factors and dimensionality reduction

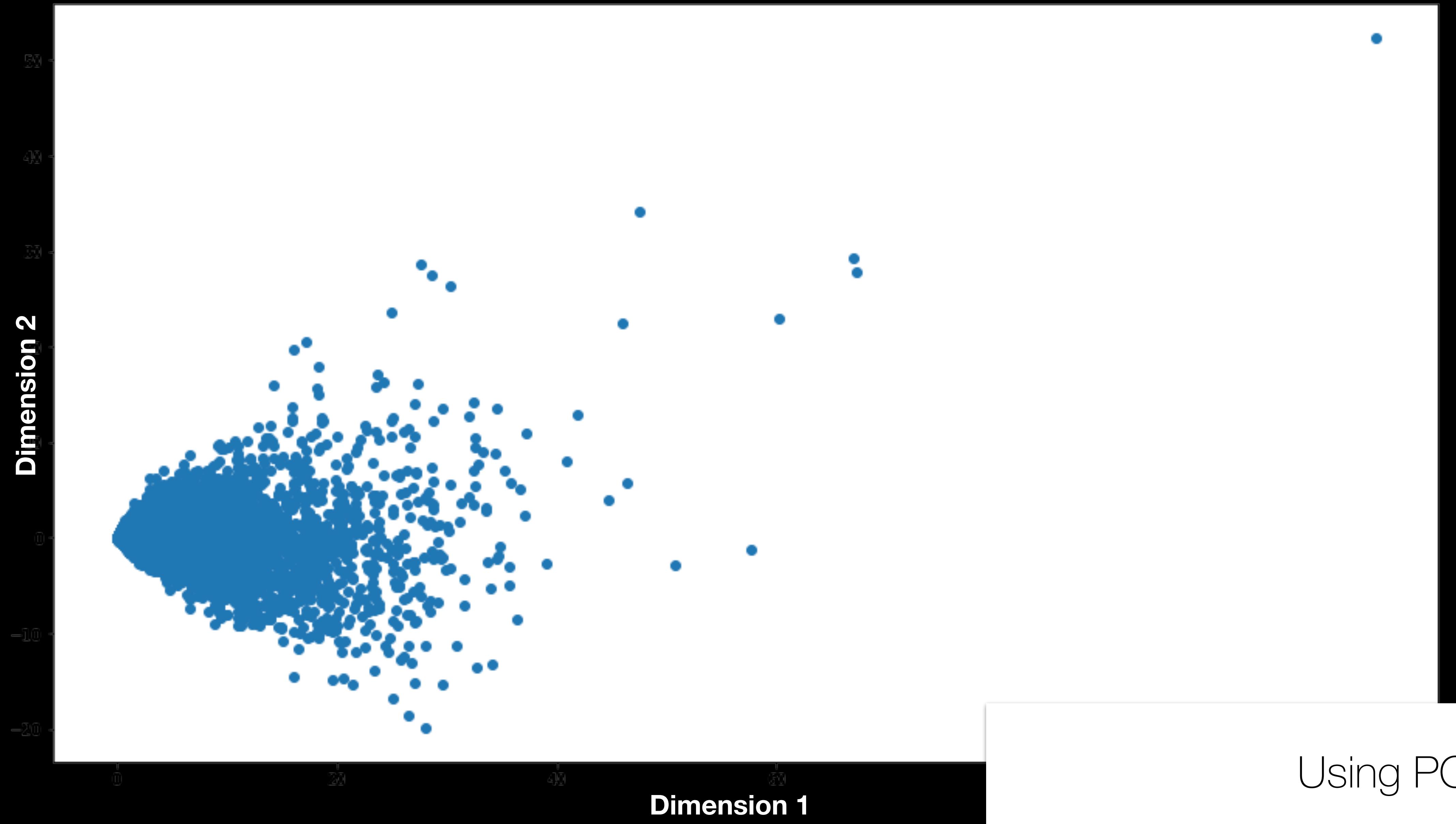
Describe at least two methods for dimensionality reduction



actor\_genre\_df.head(20)

	Comedy	Fantasy	Romance	Action	Crime	Adventure	Mystery	Thriller	Drama	Biography	...	Sport	News	Family	Western	Short
nm0000212	16.0	3.0	16.0	5.0	4.0	2.0	5.0	3.0	16.0	2.0	...	0.0	0.0	0.0	0.0	0.0
nm0413168	8.0	3.0	6.0	14.0	6.0	11.0	5.0	2.0	13.0	5.0	...	0.0	0.0	0.0	0.0	0.0
nm0000630	10.0	2.0	6.0	4.0	1.0	2.0	2.0	4.0	17.0	6.0	...	4.0	1.0	1.0	0.0	0.0
nm0005227	12.0	1.0	3.0	2.0	0.0	3.0	0.0	1.0	5.0	1.0	...	1.0	0.0	0.0	0.0	0.0
nm0697338	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm1300519	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0940707	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0625977	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0792032	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0496571	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm2868805	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm2866192	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0001379	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2.0	1.0	...	1.0	0.0	0.0	1.0	0.0
nm0462648	1.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	2.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0000953	6.0	0.0	0.0	1.0	3.0	0.0	0.0	2.0	9.0	7.0	...	0.0	0.0	0.0	0.0	0.0
nm0001782	0.0	0.0	0.0	1.0	0.0	1.0	0.0	1.0	1.0	1.0	...	0.0	0.0	0.0	0.0	0.0
nm0005077	1.0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	1.0	0.0
nm0550626	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0177016	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
nm0907480	0.0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0

How might you visualize  
these rows?



Using PCA/SVD

Development/github/umd.inst414/Module03/

scikit-learn.org/stable/modules/classes.html#module-sklearn.decomp

04-Dimensionality.PCA-Copy1 - Jupyter Notebook

API Reference — scikit-learn 1.2.1 documentation

## sklearn.decomposition: Matrix Decomposition

The `sklearn.decomposition` module includes matrix decomposition algorithms, including among others PCA, NMF or ICA. Most of the algorithms of this module can be regarded as dimensionality reduction techniques.

**User guide:** See the [Decomposing signals in components \(matrix factorization problems\)](#) section for further details.

<code>decomposition.DictionaryLearning([...])</code>	Dictionary learning.
<code>decomposition.FactorAnalysis([n_components, ...])</code>	Factor Analysis (FA).
<code>decomposition.FastICA([n_components, ...])</code>	FastICA: a fast algorithm for Independent Component Analysis.
<code>decomposition.IncrementalPCA([n_components, ...])</code>	Incremental principal components analysis (IPCA).
<code>decomposition.KernelPCA([n_components, ...])</code>	Kernel Principal component analysis (KPCA) [R396fc7d924b8-1].
<code>decomposition.LatentDirichletAllocation([...])</code>	Latent Dirichlet Allocation with online variational Bayes algorithm.
<code>decomposition.MiniBatchDictionaryLearning([...])</code>	Mini-batch dictionary learning.
<code>decomposition.MiniBatchSparsePCA([...])</code>	Mini-batch Sparse Principal Components Analysis.
<code>decomposition.NMF([n_components, init, ...])</code>	Non-Negative Matrix Factorization (NMF).
<code>decomposition.MiniBatchNMF([n_components, ...])</code>	Mini-Batch Non-Negative Matrix Factorization (NMF).
<code>decomposition.PCA([n_components, copy, ...])</code>	Principal component analysis (PCA).
<code>decomposition.SparsePCA([n_components, ...])</code>	Sparse Principal Components Analysis (SparsePCA).
<code>decomposition.SparseCoder(dictionary, *[...])</code>	Sparse coding.
<code>decomposition.TruncatedSVD([n_components, ...])</code>	Dimensionality reduction using truncated SVD (aka LSA).
<code>decomposition.dict_learning(X, n_components, ...)</code>	Solve a dictionary learning matrix factorization problem.
<code>decomposition.dict_learning_online(X[, ...])</code>	Solve a dictionary learning matrix factorization problem online.
<code>decomposition.fastica(X[, n_components, ...])</code>	Perform Fast Independent Component Analysis.
<code>decomposition.non_negative_factorization(X)</code>	Compute Non-negative Matrix Factorization (NMF).
<code>decomposition.sparse_encode(X, dictionary, *)</code>	Sparse coding.

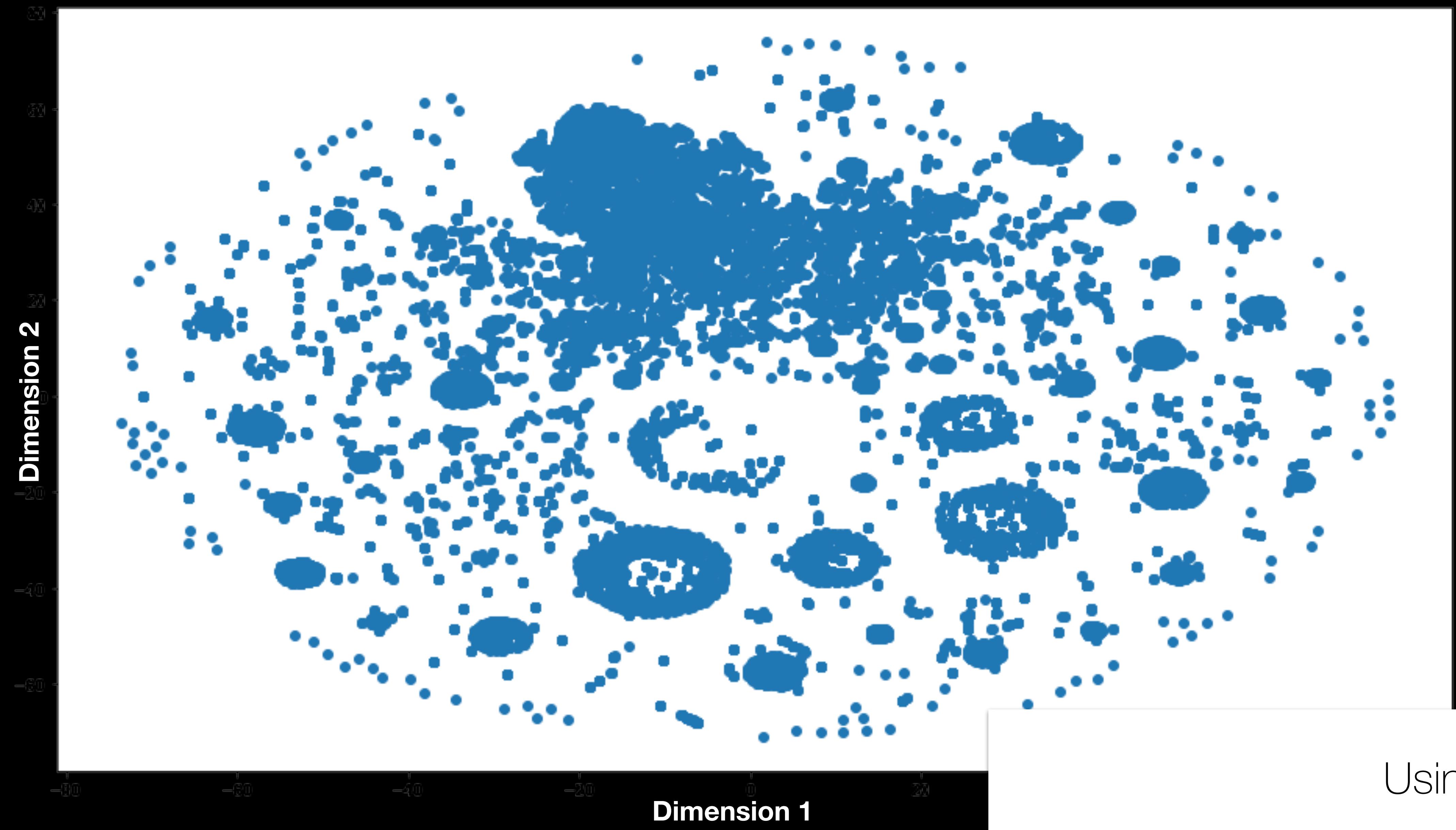
## sklearn.discriminant\_analysis: Discriminant Analysis

Linear Discriminant Analysis and Quadratic Discriminant Analysis

**User guide:** See the [Linear and Quadratic Discriminant Analysis](#) section for further details.

<code>discriminant_analysis.LinearDiscriminantAnalysis([...])</code>	Linear Discriminant Analysis.
<code>discriminant_analysis.QuadraticDiscriminantAnalysis(*)</code>	Quadratic Discriminant Analysis.

## sklearn.dummy: Dummy estimators



Using tSNE

# Questions?

Prof. Cody Buntain | @codybuntain | [cbuntain@umd.edu](mailto:cbuntain@umd.edu)  
Director, Information Ecosystems Lab