

PREVIEW OF TEXT ANALYSIS

• *Jesse Lecy* •

TEXT AS DATA

1. PRE-PROCESSING
2. TOKENIZATION
3. FEATURE SELECTION
4. MODELING

THE CORPORATION'S SPECIFIC PURPOSE IS TO SUPPORTS
AFFORDABLE HOUSING, COMMUNITY DEVELOPMENT AND
ECONOMIC DEVELOPMENT OF THE CITY AND COUNTY OF SAN
FRANCISCO'S ECONOMICALLY DISADVANTAGED INDIVIDUALS AND
COMMUNITIES, BY LENDING TO, INVESTING IN, AND DIRECTLY
ACQUIRING SUCH AFFORDABLE HOUSING AND RELATED COMMUNITY
DEVELOPMENT REAL ESTATE ASSETS.

~~the corporation specific purpose is to support~~ AFFORDABLE_HOUSING,
community development ~~and~~ ECONOMIC_DEVELOPMENT ~~of the city and county~~
of SAN_FRANCISCO economically disadvantaged individuals and communities by
lending ~~to~~ investing ~~in and~~ directly acquiring ~~such~~ AFFORDABLE_HOUSING ~~and~~
related community development REAL_ESTATE assets

1. Remove punctuation
2. Delete words with little information value (“stop words” in quanteda)
3. Identify compound constructs (apply “dictionary”)

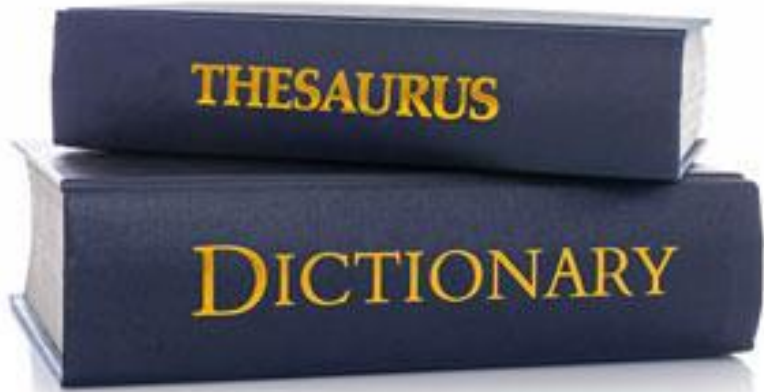
STEMMING

LEND

RELATE

LENDing

RELATED



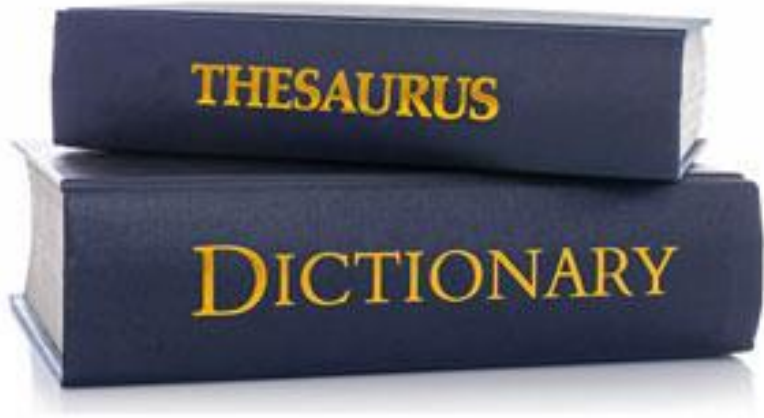
DISAMBIGUATION

George W. Bush

George Bush Jr.

President Bush

} GW_BUSH



DISAMBIGUATION



DOCUMENT FREQUENCY MATRIX (DFM): final output of pre-processing steps in quanteda

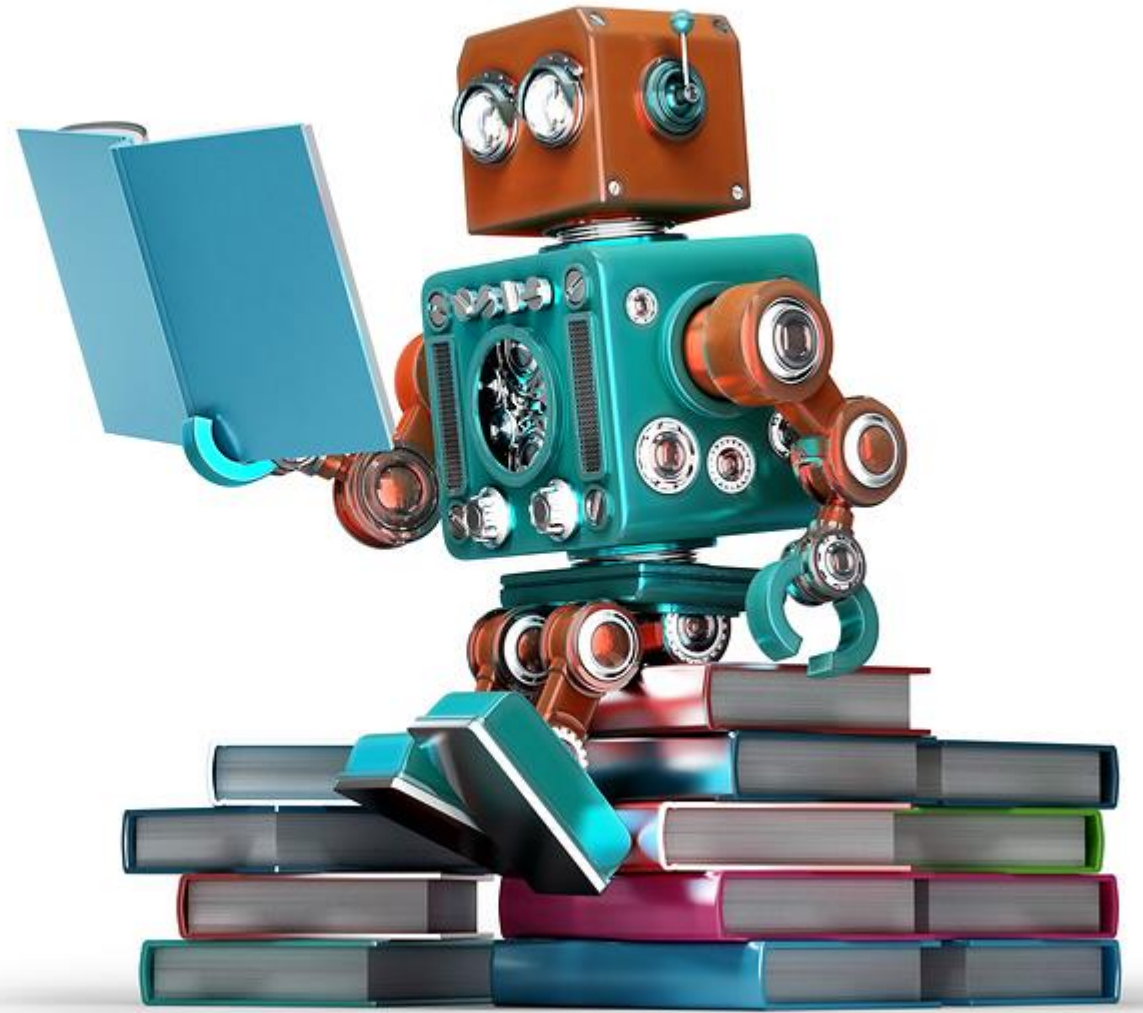
Terms	Documents													
	M	M	M	M	M	M	M	M	M	M	M	M	M	M
	1	2	3	4	5	6	7	8	9	10	11	12	13	14
abnormalities	0	0	0	0	0	0	0	1	0	1	0	0	0	0
age	1	0	0	0	0	0	0	0	0	0	0	1	0	0
behavior	0	0	0	0	1	1	0	0	0	0	0	0	0	0
blood	0	0	0	0	0	0	0	1	0	0	1	0	0	0
close	0	0	0	0	0	0	1	0	0	0	1	0	0	0
culture	1	1	0	0	0	0	0	1	1	0	0	0	0	0
depressed	1	0	1	1	1	0	0	0	0	0	0	0	0	0
discharge	1	1	0	0	0	1	0	0	0	0	0	0	0	0
disease	0	0	0	0	0	0	0	0	1	0	1	0	0	0
fast	0	0	0	0	0	0	0	0	0	1	0	1	1	1
generation	0	0	0	0	0	0	0	0	1	0	0	0	1	0
oestrogen	0	0	1	1	0	0	0	0	0	0	0	0	0	0
patients	1	1	0	1	0	0	0	1	0	0	0	0	0	0
pressure	0	0	0	0	0	0	0	0	0	0	1	0	0	1
rats	0	0	0	0	0	0	0	0	0	0	0	0	1	1
respect	0	0	0	0	0	0	0	1	0	0	0	1	0	0
rise	0	0	0	1	0	0	0	0	0	0	0	0	0	1
study	1	0	1	0	0	0	0	0	1	0	0	0	0	0

MACHINE LEARNING

(CLASSIFICATION)

TEXT AS DATA:

Text is a useful input for many machine learning models for prediction and trend analysis





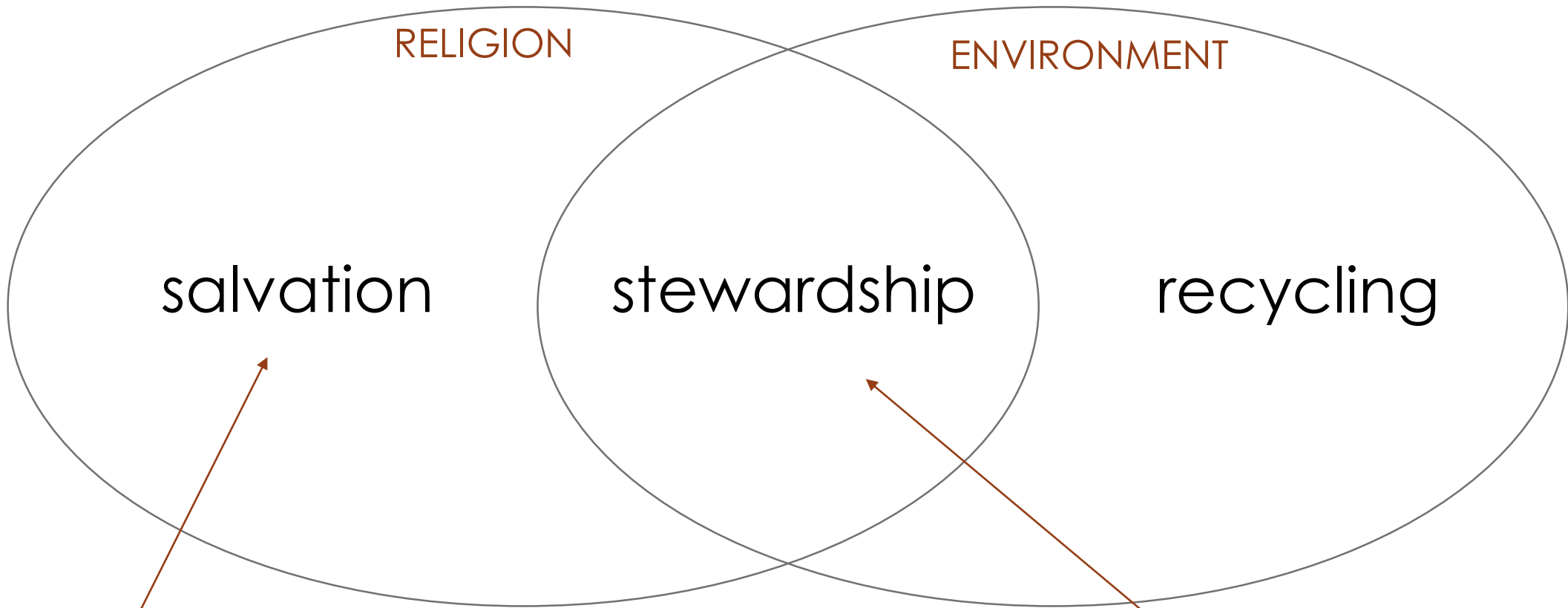
The Netflix Prize

Netflix provided a *training data set* of 100,480,507 ratings that 480,189 users gave to 17,770 movies, graded from 1 to 5 stars.



The *test set* of 1,408,789 ratings is used by the jury to determine potential prize winners.

In 2009 the *grand prize of \$1,000,000* was given to the BellKor's Pragmatic Chaos team which improved Netflix's own algorithm by 10%



Distinctive
constructs

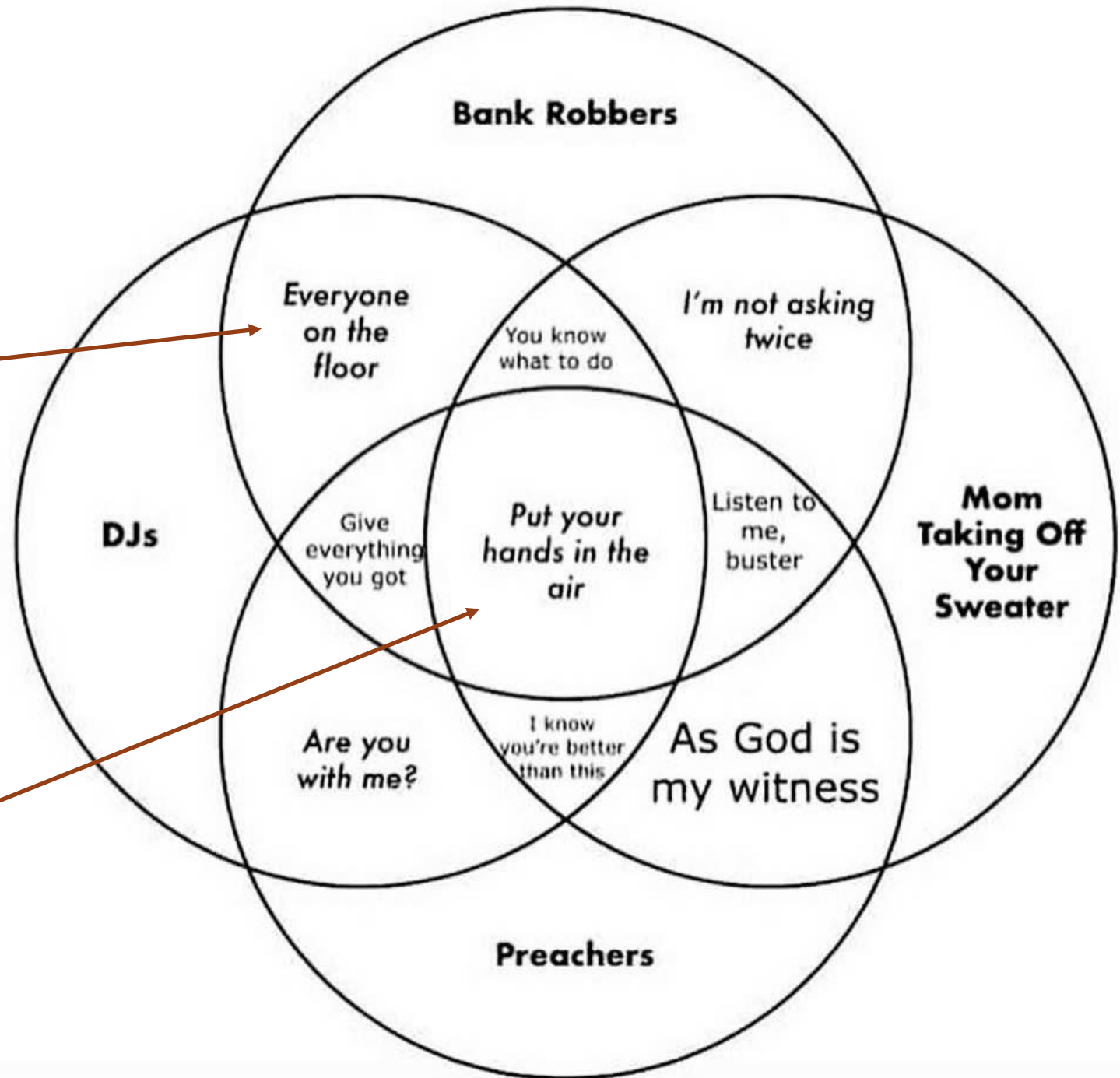
- I. Arts, Culture, and Humanities
- II. Education
- III. Environment and Animals
- IV. Health
- V. Human Services
- VI. International, Foreign Affairs
- VII. Public, Societal Benefit
- VIII. Religion Related

Shared
constructs

Information theory

HIGH Signal
to Noise Ratio
(eliminates half of
possible sources)

LOW Signal
to Noise Ratio
(does not distinguish
any actors)

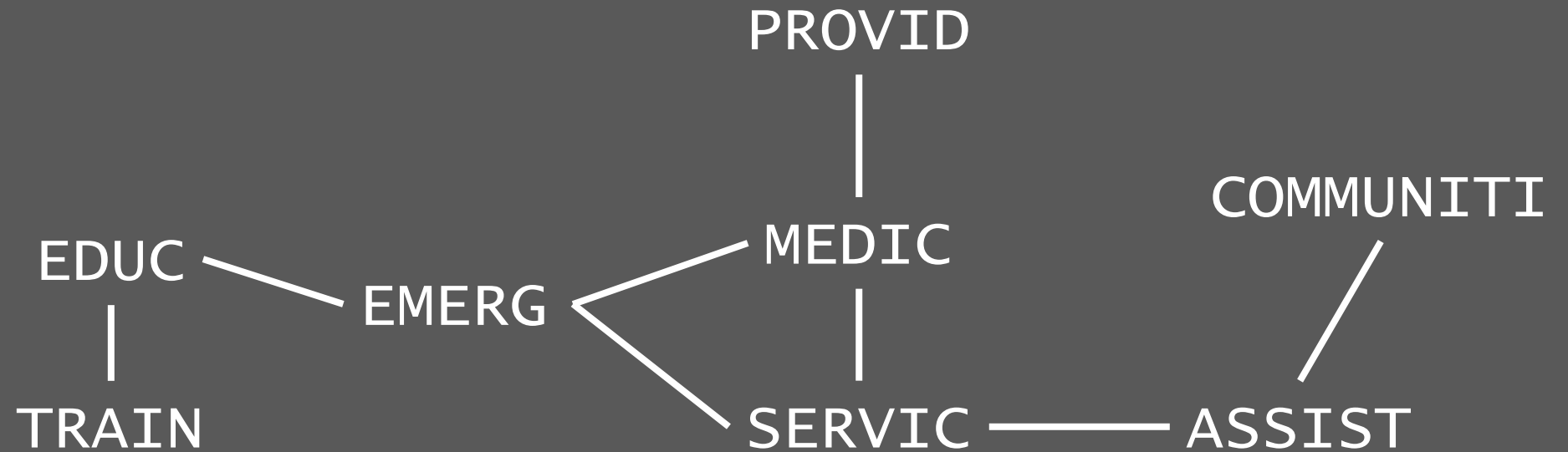


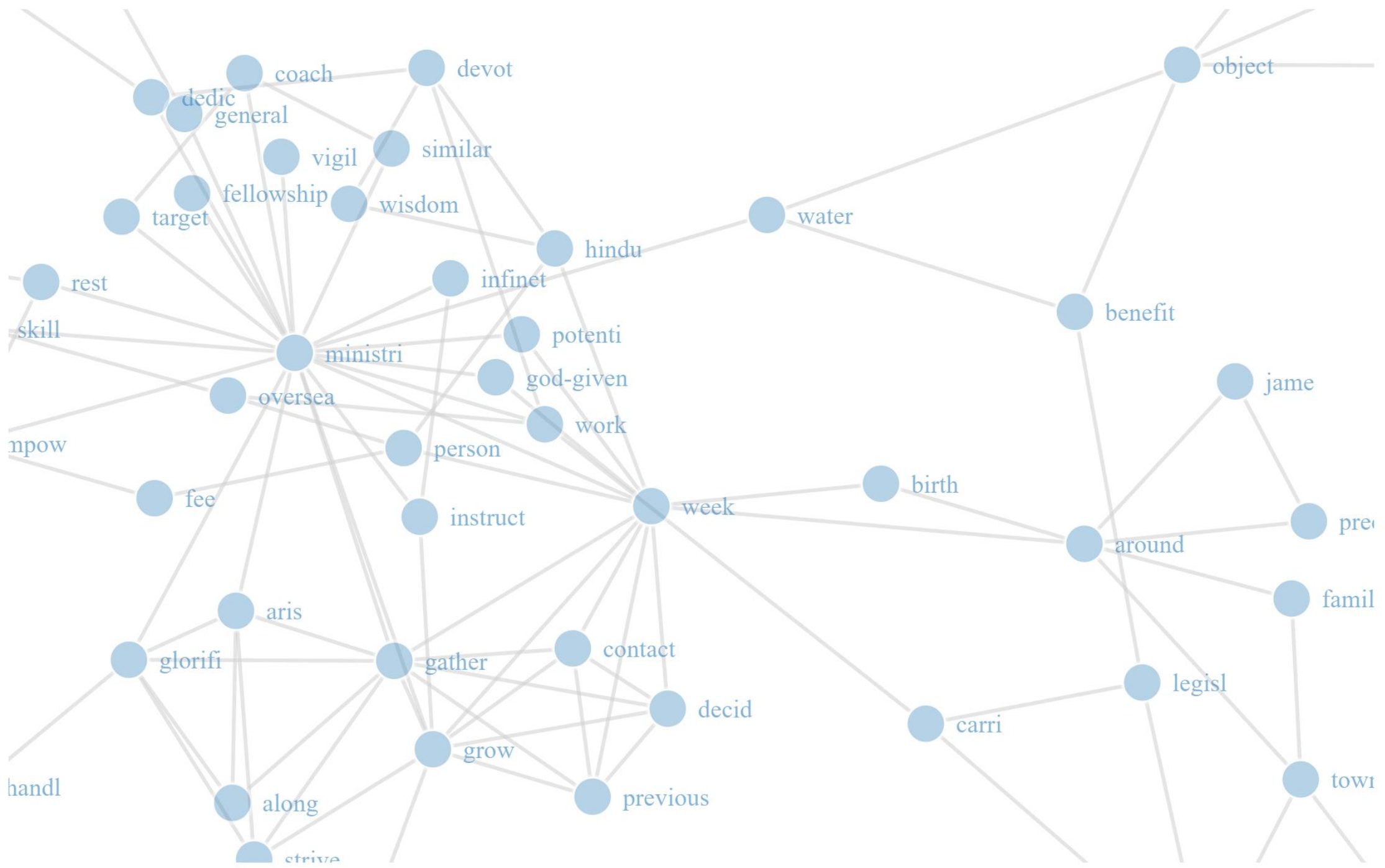
To educate, train and assist in
providing emergency medical
service for the community.



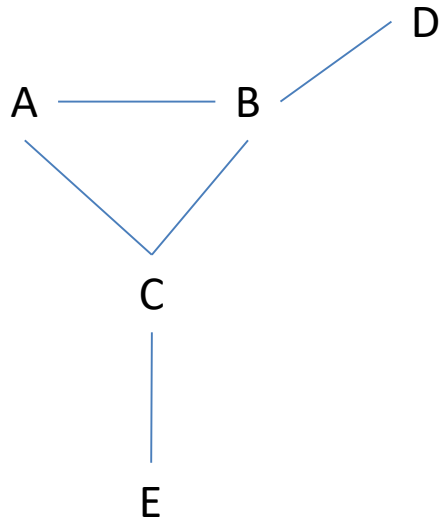
"EDUC" "TRAIN" "ASSIST"
"PROVID" "EMERG" "MEDIC"
"SERVIC" "COMMUNITI"

Semantic Networks

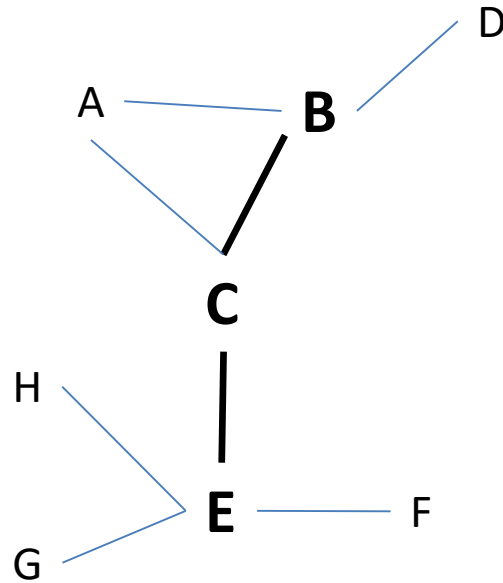




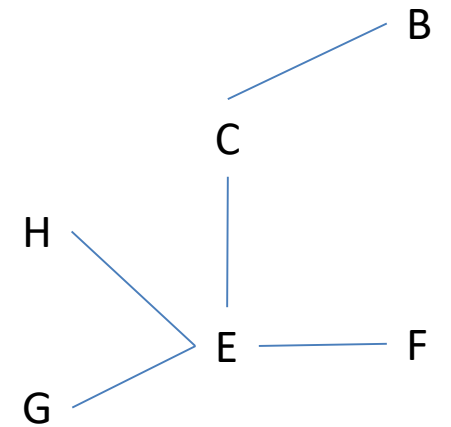
Mission Statement 1



Union (all statement) and **Intersection**



Mission Statement 2

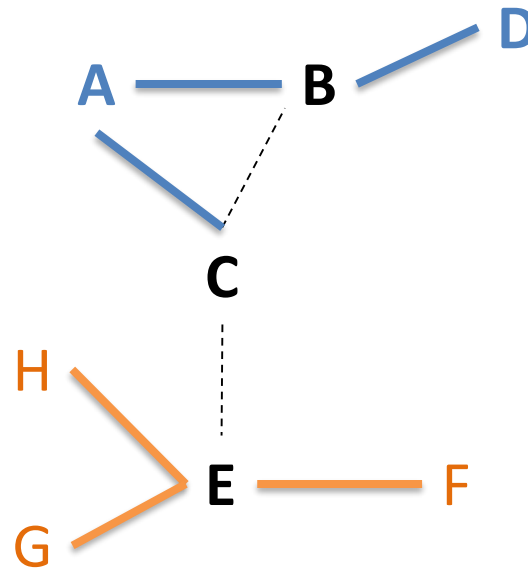


Intersection represents components of mission statements common to all nonprofits.

Analyzing Missions by Types of Nonprofits

Mission Statement
Components
Unique to Org 1:

A-B
A-C
B-D



H-E
G-E
E-F

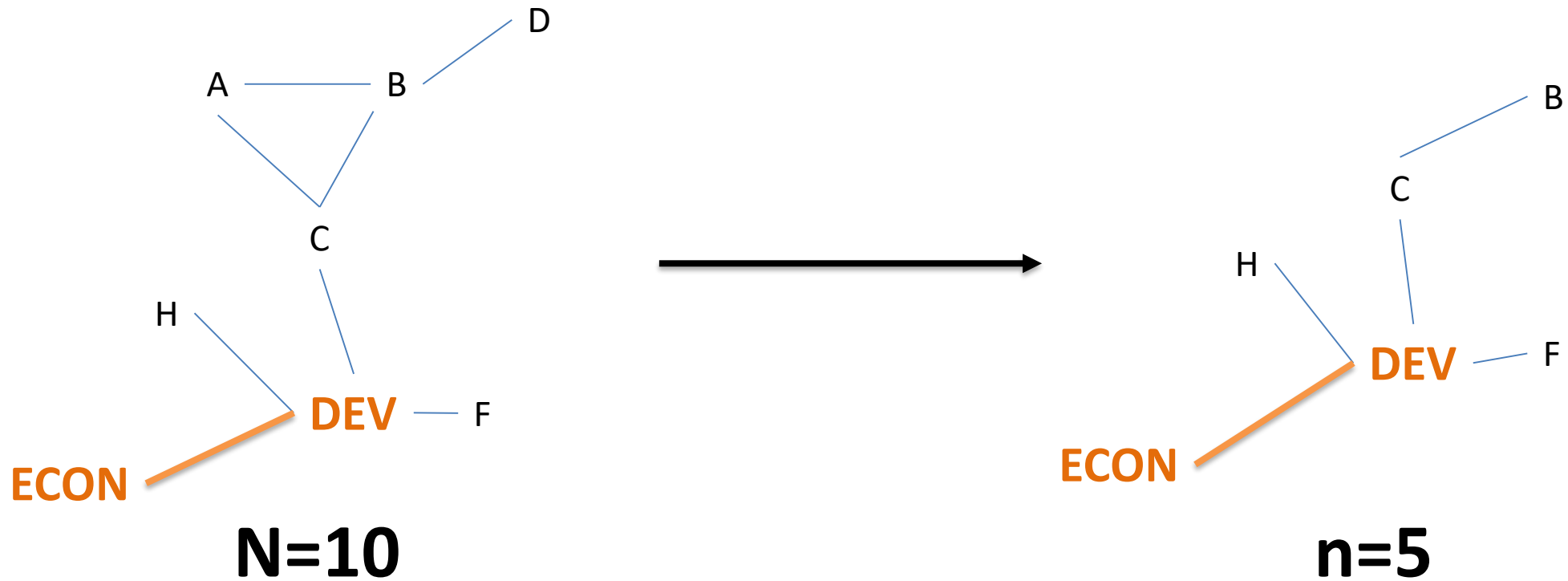
Mission Statement
Components
Unique to Org 2

Doesn't work well with dense weighted graphs!

Data structure of a weighted network:

<u>Freq ALL</u>	<u>Freq GROUP</u>	<u>Term 1</u>	<u>Term 2</u>
10	5	econ	dev
7	4	self	reliance
5	3	dev	con
5	2	globla	econ
4	2	local	econ
4	1	soc	econ
3	2	econ	socialism
3	3	finance	global
3	2	global	finance
3	2	global	impsm
3	1	impsm	global
3	1	impsm	invasion

Is it significant that **economic development** was mentioned
5 times by a specific type of organization?

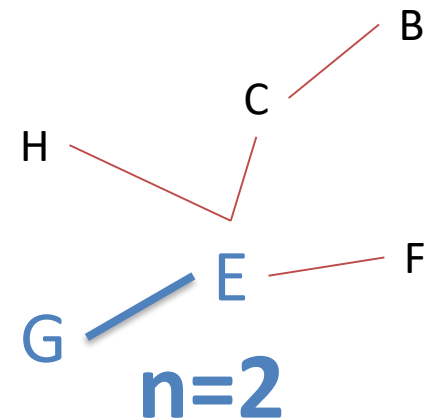
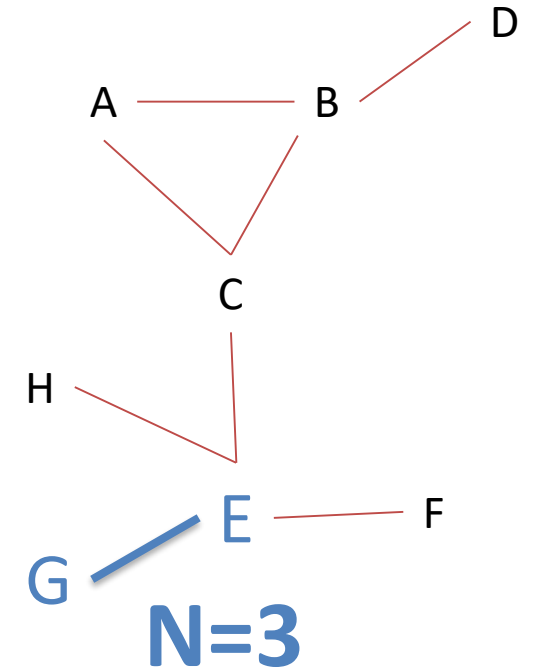
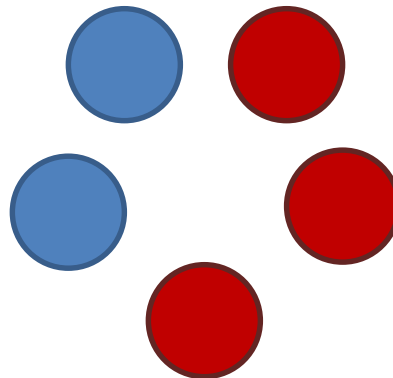
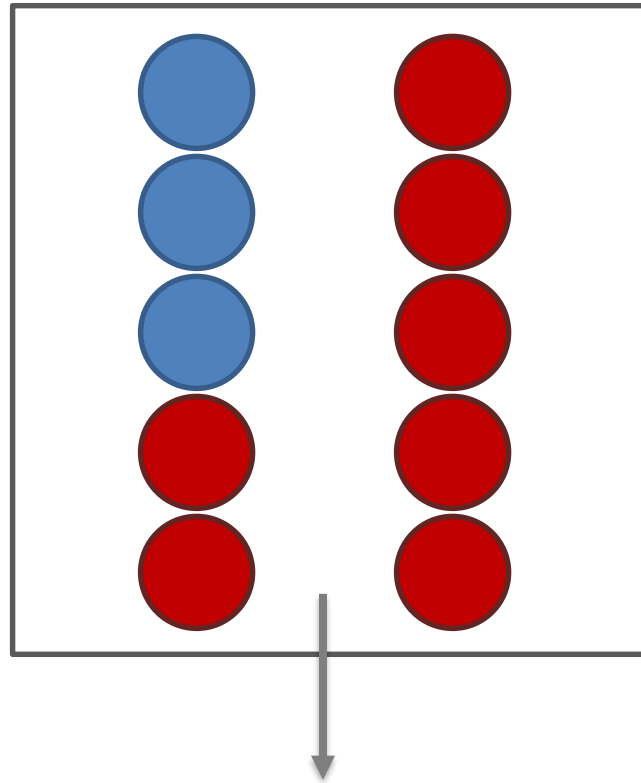


How often will a **random sample** from the full weighted network produce the **observed number** of “statements” (semantic network ties)?

Is it significant that **Org**
Statement significant?

What is the probability of
selecting **2 blue balls** from a
sample of 5 balls?

$$\Pr(\text{blue} = 2 \mid n = 5) = \frac{\binom{3}{2} \binom{7}{3}}{\binom{10}{5}} \\ = 0.42$$



Generalized:

$$\Pr(\text{StatementCount} = x \mid \text{sample} = k) = \frac{\binom{X}{x} \binom{N-X}{k-x}}{\binom{N}{k}}$$

Where X = the number of time a statement appears

N = total number of statements

k = number of statements in a specific period or group

