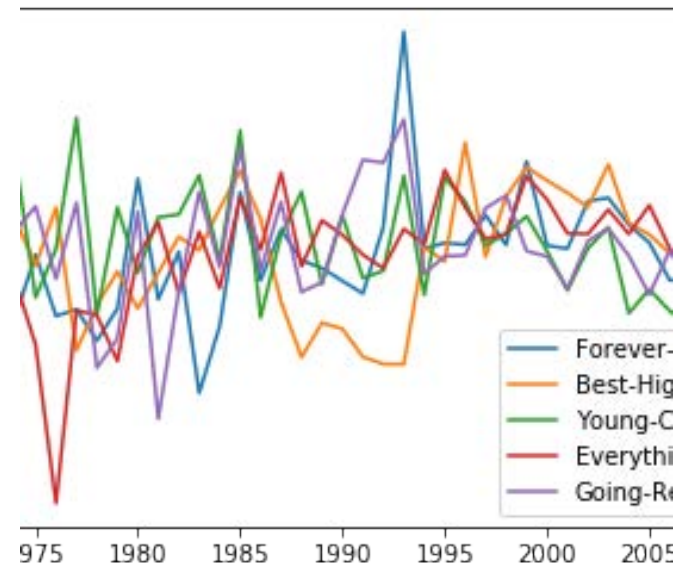# Text Intro

School of Information Studies
Syracuse University

# Finding Meaning in Text Data

- **Obtain**

- **Scrub**

- **Explore**

- Model

- iNterpret

- One of the Funniest Police Shows Ever      (1974)

- Mosaiken und Fresken als Zeugen      (1975)

- Let all the poisons that lurk in the mud      (1976)

...

- 12 days that shook Chile      (2011)

- The Movie Unrated and Uncensored      (2012)

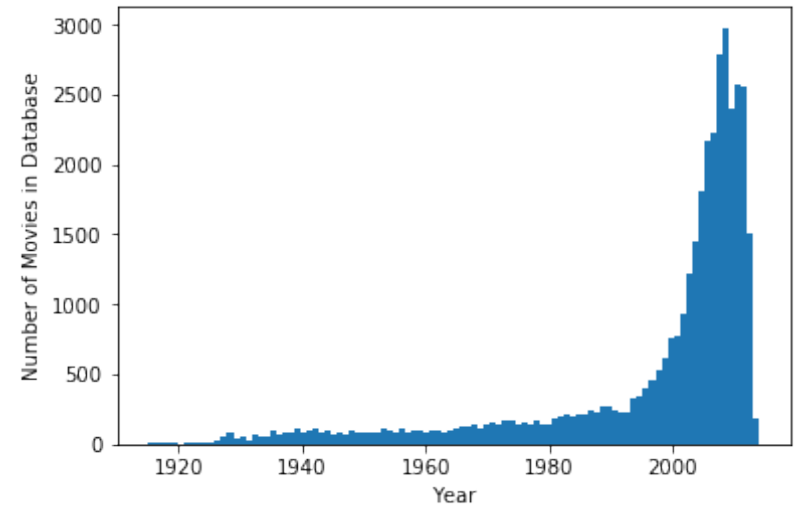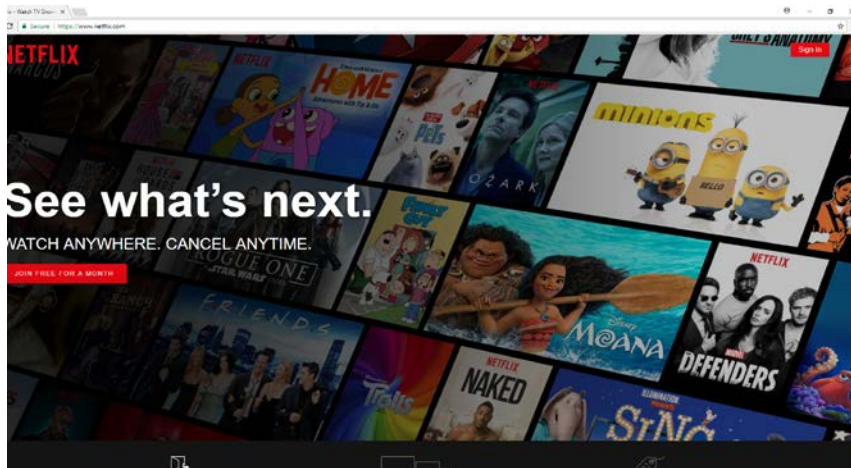- Pretty on the outside, ugly crying on the inside      (2013)

School of Information Studies
Syracuse University

# Our Challenge This Week?

School of Information Studies
Syracuse University

# What Are the Movie Trends?

School of Information Studies
Syracuse University

# Can We Use Text Analytics to Detect Trends?

School of Information Studies
Syracuse University

Data Review

School of Information Studies
Syracuse University

# Text Data



| | | | | |
|---|---|---|---|---|
| actors | often | far | years | idea |
| black | one | film | acting | kids |
| early | part | funny | action | line |
| enjoy | perhaps | goes | actually | look |
| films | say | however | back | many |
| full | supposed | long | can | movie |
| get | thought | love | characters | never |
| good | three | make | done | new |
| great | two | picture | enjoyed | night |
| guys | along | plot | especially | now |
| humor | american | quite | every | original |
| just | another | rather | excellent | people |
| last | best | script | face | plays |
| later | comedy | still | feel | read |
| lead | definitely | story | finally | real |
| like | despite | stupid | first | right |
| little | direction | war | gives | role |
| loved | director | way | got | scene |
| men | end | well | heard | scenes |
| minutes | entertaining | worth | history | screen |

School of Information Studies
Syracuse University

# Bag of Words

things party
know mind
hell
wild *girls* beyond
stop past *high* *going* big
*want*
*young*
great men
family get young
everything year movie
game man
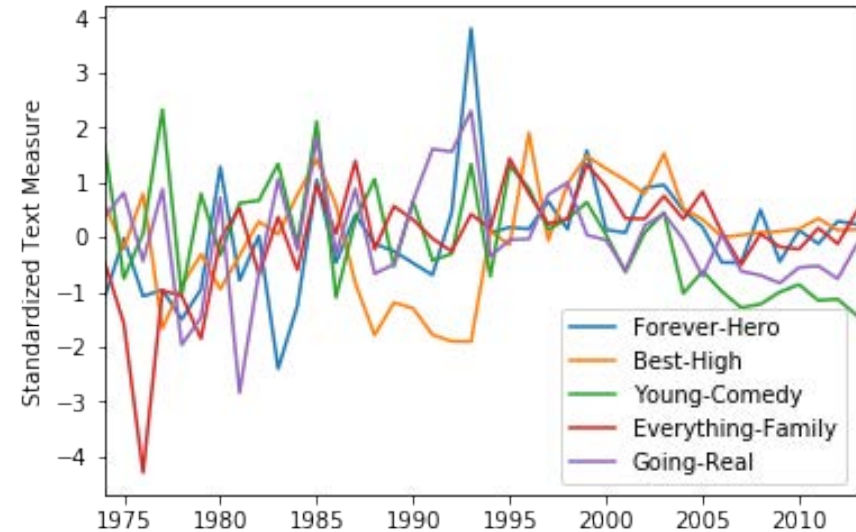
- One of the Funniest Police Shows Ever          (1974)
- Mosaiken und Fresken als Zeugen          (1975)
- Let all the poisons that lurk in the mud          (1976)

...

- 12 Days That Shook Chile          (2011)
- The Movie Unrated and Uncensored          (2012)
- Pretty on the Outside, Ugly Crying on the Inside     (2013)

School of Information Studies
Syracuse University

# Bag of Words (cont.)

School of Information Studies
Syracuse University

# Sentiment Analysis

| | | | | |
|---|---|---|---|---|
| actors | often | far | years | idea |
| black | one | film | acting | kids |
| early | part | funny | action | line |
| enjoy | perhaps | goes | actually | look |
| films | say | however | back | many |
| full | supposed | long | can | movie |
| get | thought | love | characters | never |
| good | three | make | done | new |
| great | two | picture | enjoyed | night |
| guys | along | plot | especially | now |
| humor | american | quite | every | original |
| just | another | rather | excellent | people |
| last | best | script | face | plays |
| later | comedy | still | feel | read |
| lead | definitely | story | finally | real |
| like | despite | stupid | first | right |
| little | direction | war | gives | role |
| loved | director | way | got | scene |
| men | end | well | heard | scenes |
| minutes | entertaining | worth | history | screen |

School of Information Studies
Syracuse University

Recommendation | School of Information Studies
Syracuse University

# Our Challenge This Week?

**School of Information Studies**
Syracuse University

# How Might We Discover Them?

School of Information Studies
Syracuse University

# How Might We Discover Them? (cont.)

School of Information Studies
Syracuse University

# Bag of Words

| | | | | |
|---|---|---|---|---|
| actors | often | far | years | idea |
| black | one | film | acting | kids |
| early | part | funny | action | line |
| enjoy | perhaps | goes | actually | look |
| films | say | however | back | many |
| full | supposed | long | can | movie |
| get | thought | love | characters | never |
| good | three | make | done | new |
| great | two | picture | enjoyed | night |
| guys | along | plot | especially | now |
| humor | american | quite | every | original |
| just | another | rather | excellent | people |
| last | best | script | face | plays |
| later | comedy | still | feel | read |
| lead | definitely | story | finally | real |
| like | despite | stupid | first | right |
| little | direction | war | gives | role |
| loved | director | way | got | scene |
| men | end | well | heard | scenes |
| minutes | entertaining | worth | history | screen |

things party
wild
stop past

high girls going
want
young

know mind
hell
beyond
big

family great
everything year
game

men young
get
movie
man

School of Information Studies
Syracuse University

# Sentiment Analysis

| | | | | |
|---|---|---|---|---|
| actors | often | far | years | idea |
| black | one | film | acting | kids |
| early | part | funny | action | line |
| enjoy | perhaps | goes | actually | look |
| films | say | however | back | many |
| full | supposed | long | can | movie |
| get | thought | love | characters | never |
| good | three | make | done | new |
| great | two | picture | enjoyed | night |
| guys | along | plot | especially | now |
| humor | american | quite | every | original |
| just | another | rather | excellent | people |
| last | best | script | face | plays |
| later | comedy | still | feel | read |
| lead | definitely | story | finally | real |
| like | despite | stupid | first | right |
| little | direction | war | gives | role |
| loved | director | way | got | scene |
| men | end | well | heard | scenes |
| minutes | entertaining | worth | history | screen |

things  party
wild
stop  past

girls
high  going
want
young

know  mind
hell
beyond
big

great
family
everything  year
game

men
get  young
movie
man

School of Information Studies
Syracuse University

# Recommendation

- Text analytics requires a move from unstructured to structured
- Some primary methods
  - Bag of words
  - Natural language processing
  - Sentiment analysis
- Text analytics gives new insight into consumer preferences

School of Information Studies
Syracuse University

# Text Processing

School of Information Studies
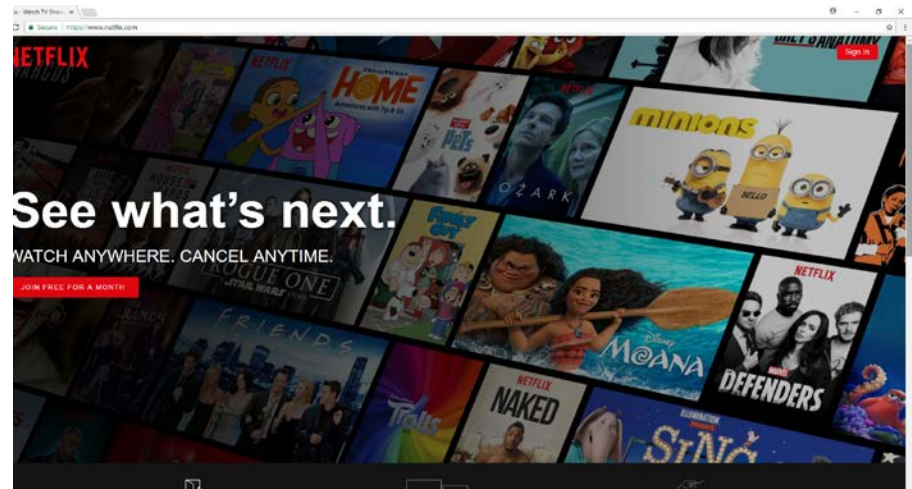Syracuse University

# Patterns in Text Data

| | | | | |
|---|---|---|---|---|
| actors | often | far | years | idea |
| black | one | film | acting | kids |
| early | part | funny | action | line |
| enjoy | perhaps | goes | actually | look |
| films | say | however | back | many |
| full | supposed | long | can | movie |
| get | thought | love | characters | never |
| good | three | make | done | new |
| great | two | picture | enjoyed | night |
| guys | along | plot | especially | now |
| humor | american | quite | every | original |
| just | another | rather | excellent | people |
| last | best | script | face | plays |
| later | comedy | still | feel | read |
| lead | definitely | story | finally | real |
| like | despite | stupid | first | right |
| little | direction | war | gives | role |
| loved | director | way | got | scene |
| men | end | well | heard | scenes |
| minutes | entertaining | worth | history | screen |

School of Information Studies
Syracuse University

# Text Processing



Source: Adapted from Miller (2005).

School of Information Studies
Syracuse University

# Term Document Matrix

Covert (2017)

You don't get more back to the future than this wholehearted embrace of sci-fi and golden oldies nostalgia. No franchise is so gifted at exploiting the popularity of superhero movies while satirizing the genre … "Guardians of the Galaxy Vol. 2" is just right. It's hyperbolic nonsense wrapped in the colors of a neon rainbow, bouncing from one artfully wacky scenario to the next. Here it's sleekly futuristic, there it's older than mud. It's the galaxy's silliest thrill show.

| Term | Document | | |
|---|---|---|---|
| | Covert | Lapin | ... |
| future | 1 | 0 | |
| oldies | 1 | 1 | |
| movie | 1 | 1 | |
| galaxy | 1 | 1 | |
| neon | 1 | 0 | |
| rainbow | 1 | 0 | |
| battle | 0 | 1 | |
| families | 0 | 1 | |
| ... | | | |

Lapin (2017)

Swapping out yet another giant, ponderous battle for a cute creature and Oldies radio is exactly the sort of action movie send-up the genre needs right now, so it's a shame the rest of the film doesn't follow that same spirit of rule-breaking. Instead, we get straight-laced morals about families that feel like warmed-over *Fast and Furious* wisdom.

Source: Adapted from Miller (2005).
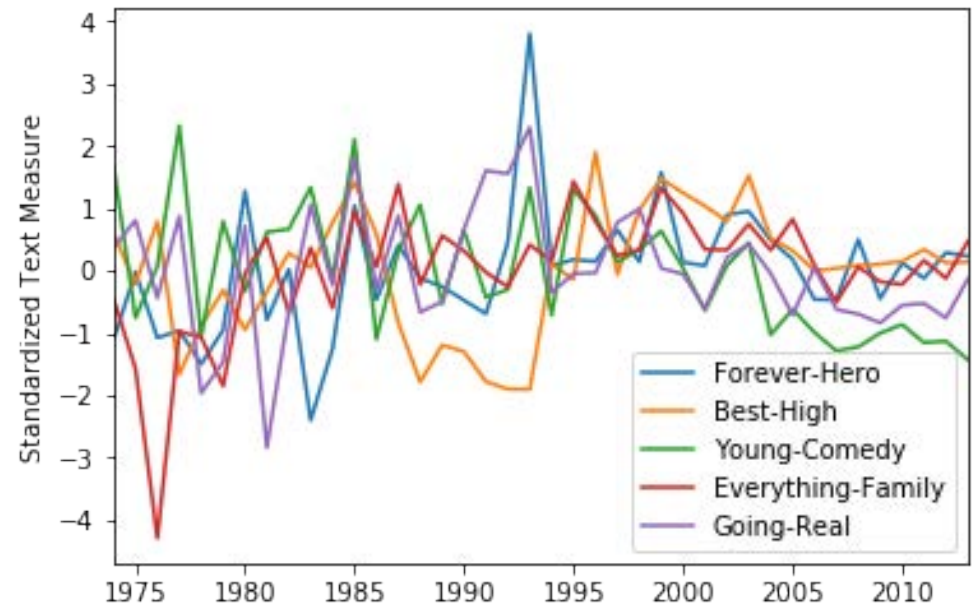
School of Information Studies
Syracuse University

# Text Analytics

Unsupervised

- No response or class
- Patterns or trends over time/population
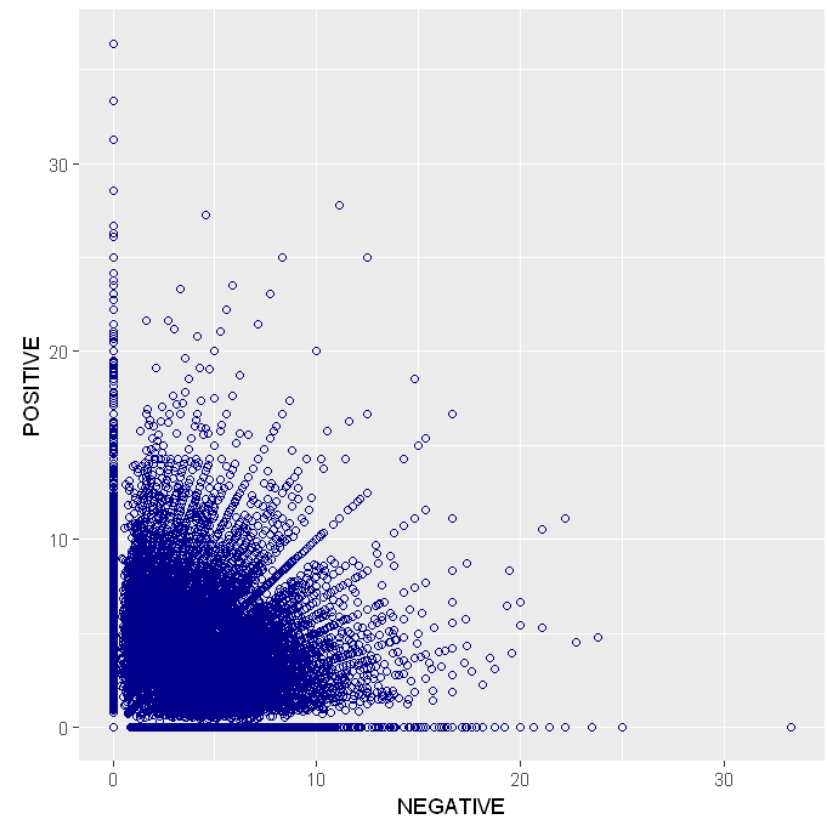
Supervised

- Response or class
- Spam or relevant search

School of Information Studies
Syracuse University

# Foundation

School of Information Studies
Syracuse University

# Patterns in Text Data

School of Information Studies
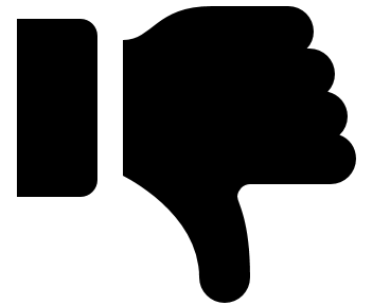Syracuse University

# Sentiment Analysis

- Measurement focused
  - Positive words
  - Negative words
- Dictionary
  - Not alphabetized definitions
  - Repository of lists
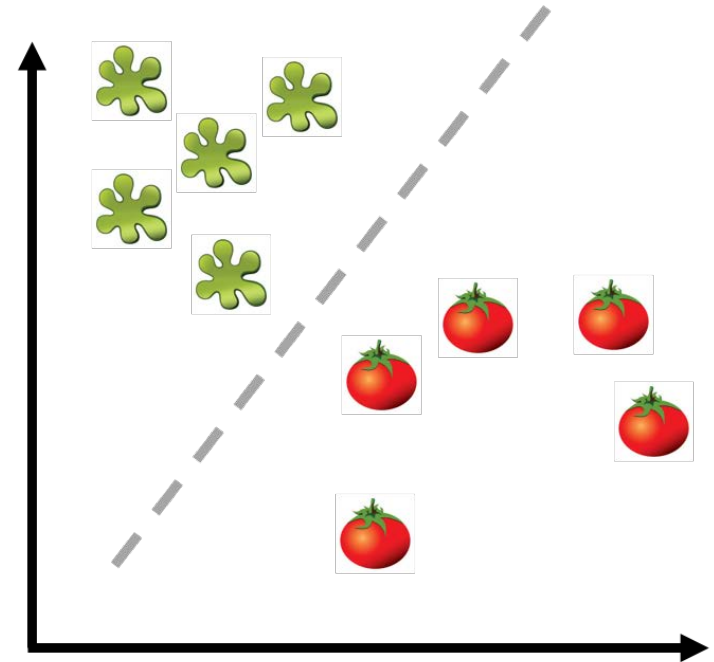  - Positive and negative words
  - Synonyms and antonyms

School of Information Studies
Syracuse University

# Simple Difference

- Compute difference scores

- Cutoff predicted on training set

- Positive minus negative scores

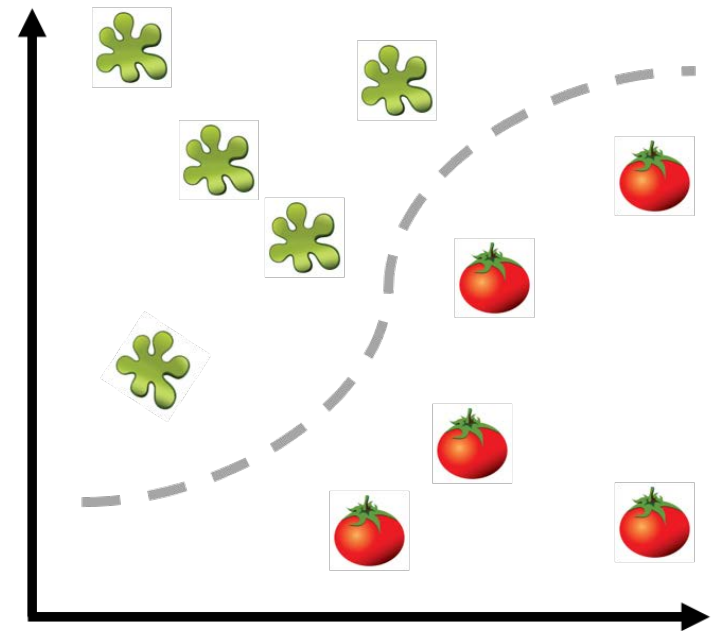School of Information Studies
Syracuse University

# Regression Difference

- Linear regression to determine weights

- Combine positive and negative scores into a linear predictor

- Predicted rating becomes a cutoff

School of Information Studies
Syracuse University
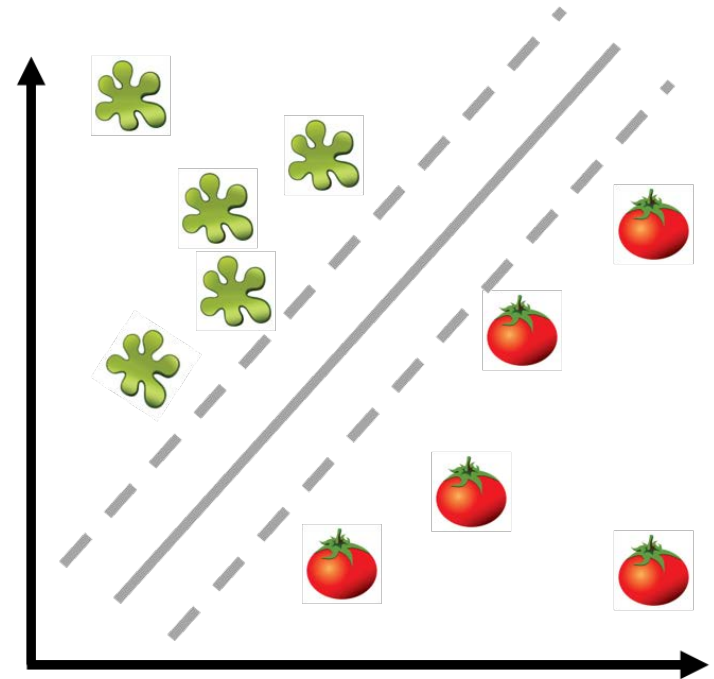
# Logistic Regression

- Predicting a binary response

- Select useful predictors from sentiment words

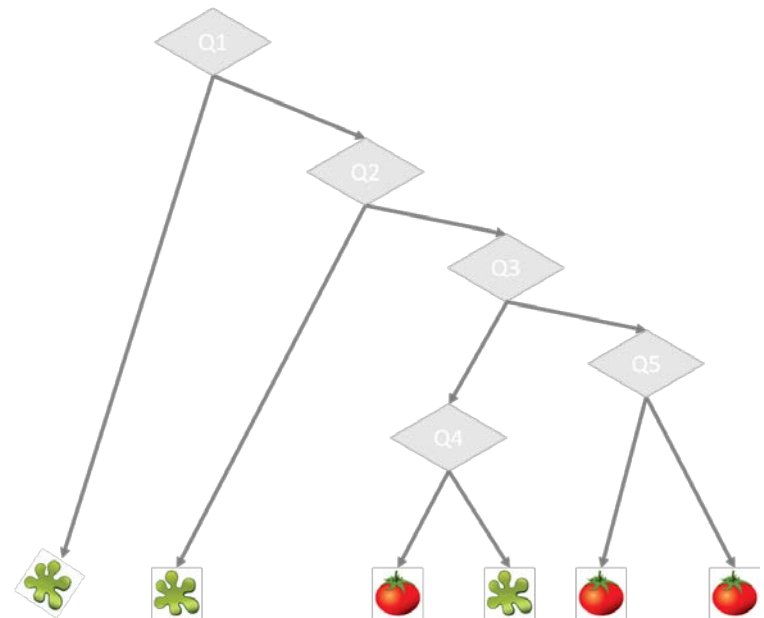- Weights calculated using maximum likelihood

# Support Vector Machines

- Effective in problems with large numbers of explanatory variables

- Classifies using a separating hyperplane

- Memory efficient

School of Information Studies
Syracuse University

# Random Forests

- Ensemble method using multiple decision trees

- Recursive partitioning on the training set

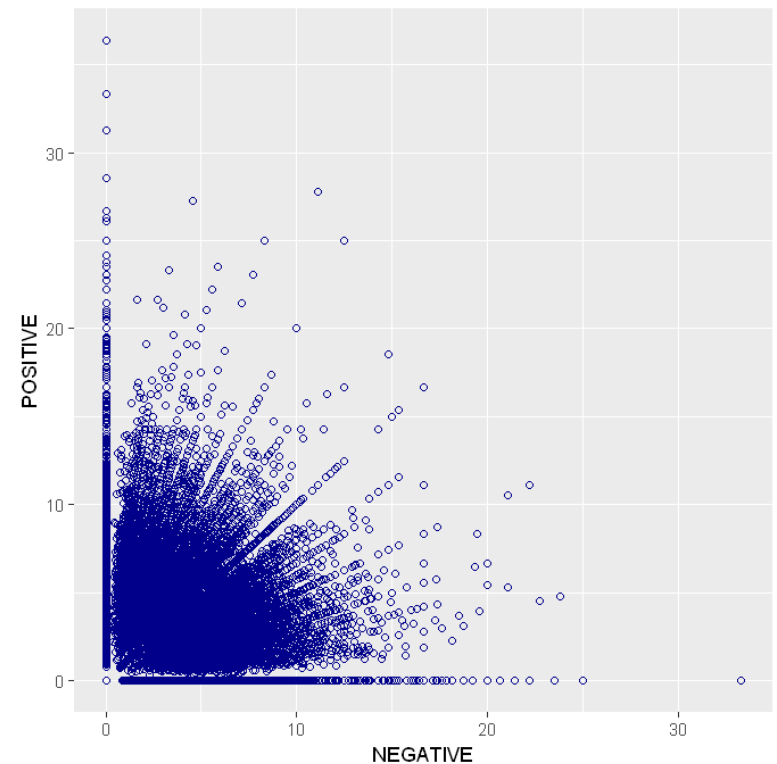- Effective with large number of explanatory variables

School of Information Studies
Syracuse University

# Alternatives
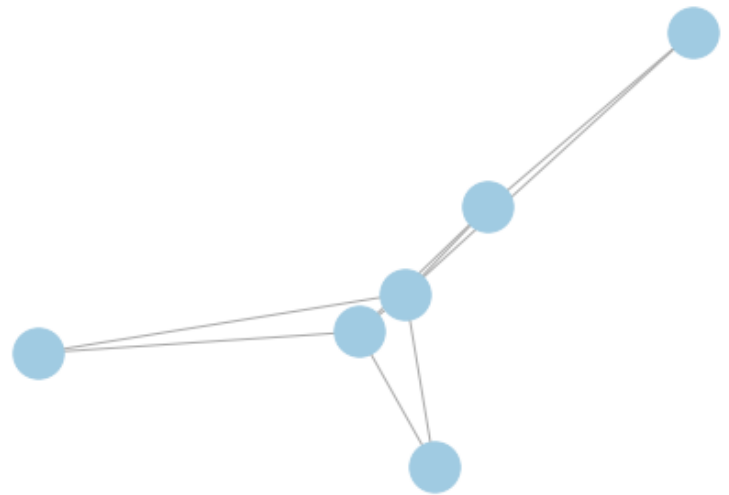
# Alternatives for Text Data

School of Information Studies
Syracuse University

# Latent Semantic Analysis

- Identify relationships between words

- Concepts within the documents

- Eliminate the noise and provide simple representation of the data

School of Information Studies
Syracuse University

# Latent Semantic Analysis (cont.)

**Willmore (2017)**

At various points during *Guardians of the Galaxy Vol. 2*, Gamora (Zoe Saldana) wields a gun the size of a tree trunk, Yondu (Michael Rooker) slaughters dozens upon dozens with his whistle-controlled arrow, and high priestess Ayesha (Elizabeth Debicki) presides over a legion of pilots controlling droids like the manager of a deadly arcade. But no weapon can compete with the small, shining eyes of Baby Groot, a creation whose immaculately crafted adorableness could level cities and sink small-to-medium continents.

**Orr (2017)**

Perhaps the finest, funniest moment in *Guardians of the Galaxy Vol. 2* is the first action sequence. Or perhaps I should put quote marks around that: "action sequence." Because for most of its duration, the action is strictly an afterthought. The titular supergroup has been enlisted to defeat a giant star-squid, and its smallest member, Baby Groot (the twig-like offshoot of last installment's arboreal giant), is hooking up some equipment in the foreground as the fight commences behind him. What is Baby Groot fiddling with? Some kind of space cannon?

**Covert (2017)**

You don't get more back to the future than this wholehearted embrace of sci-fi and golden oldies nostalgia. No franchise is so gifted at exploiting the popularity of superhero movies while satirizing the genre … "Guardians of the Galaxy Vol. 2" is just right. It's hyperbolic nonsense wrapped in the colors of a neon rainbow, bouncing from one artfully wacky scenario to the next. Here it's sleekly futuristic, there it's older than mud. It's the galaxy's silliest thrill show.

*GG2 delivers easy summer movie fare, has a message of family, stars Baby Groot, and is generally liked by reviewers.*

**Lapin (2017)**

Swapping out yet another giant, ponderous battle for a cute creature and Oldies radio is exactly the sort of action movie send-up the genre needs right now, so it's a shame the rest of the film doesn't follow that same spirit of rule-breaking. Instead, we get straight-laced morals about families that feel like warmed-over *Fast and Furious* wisdom.

# Latent Semantic Analysis (cont.)

From our term document matrix:

$$M = USV^T$$

$M$ is $m \times n$

$U$ is $m \times k$

$S$ is $k \times k$

$V$ is $k \times n$

| Term | Document | | |
|---|---|---|---|
| | Covert | Lapin | ... |
| future | 1 | 0 | |
| oldies | 1 | 1 | |
| movie | 1 | 1 | |
| galaxy | 1 | 1 | |
| neon | 1 | 0 | |
| rainbow | 1 | 0 | |
| battle | 0 | 1 | |
| families | 0 | 1 | |
| ... | | | |

School of Information Studies
Syracuse University

# Some Warnings

- The "good-bad" dimension

- Specific terms vs. latent semantic dimensions

- Relation to business outcomes and insights

School of Information Studies
Syracuse University