

TOPIC MODELING WITH STRUCTURED PRIORS FOR TEXT-DRIVEN SCIENCE

MICHAEL J. PAUL
JOHNS HOPKINS UNIVERSITY



University of Colorado, Boulder | February 27, 2015



#whiteandgold

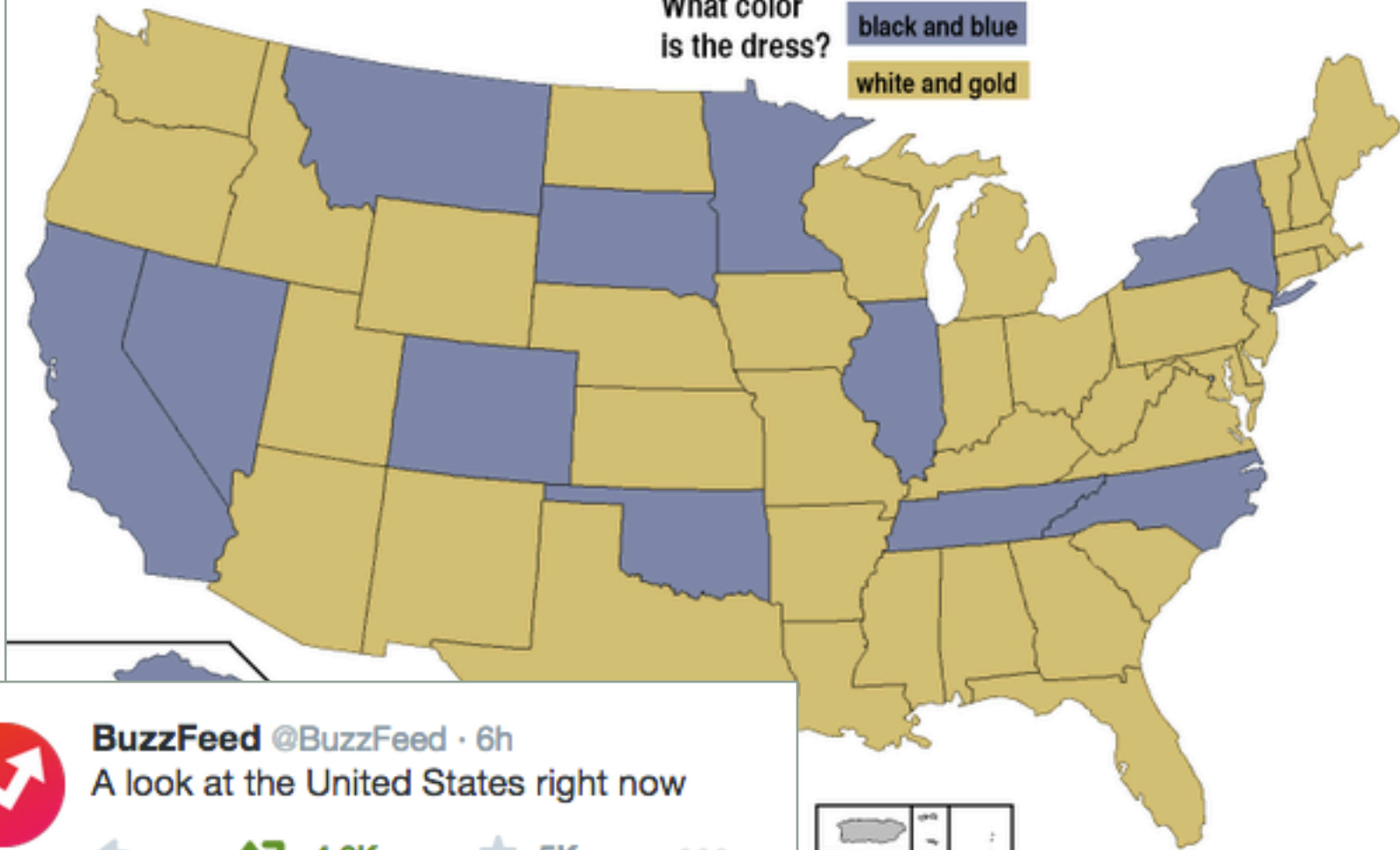


#blueandblack

What color
is the dress?

black and blue

white and gold



BuzzFeed @BuzzFeed · 6h
A look at the United States right now



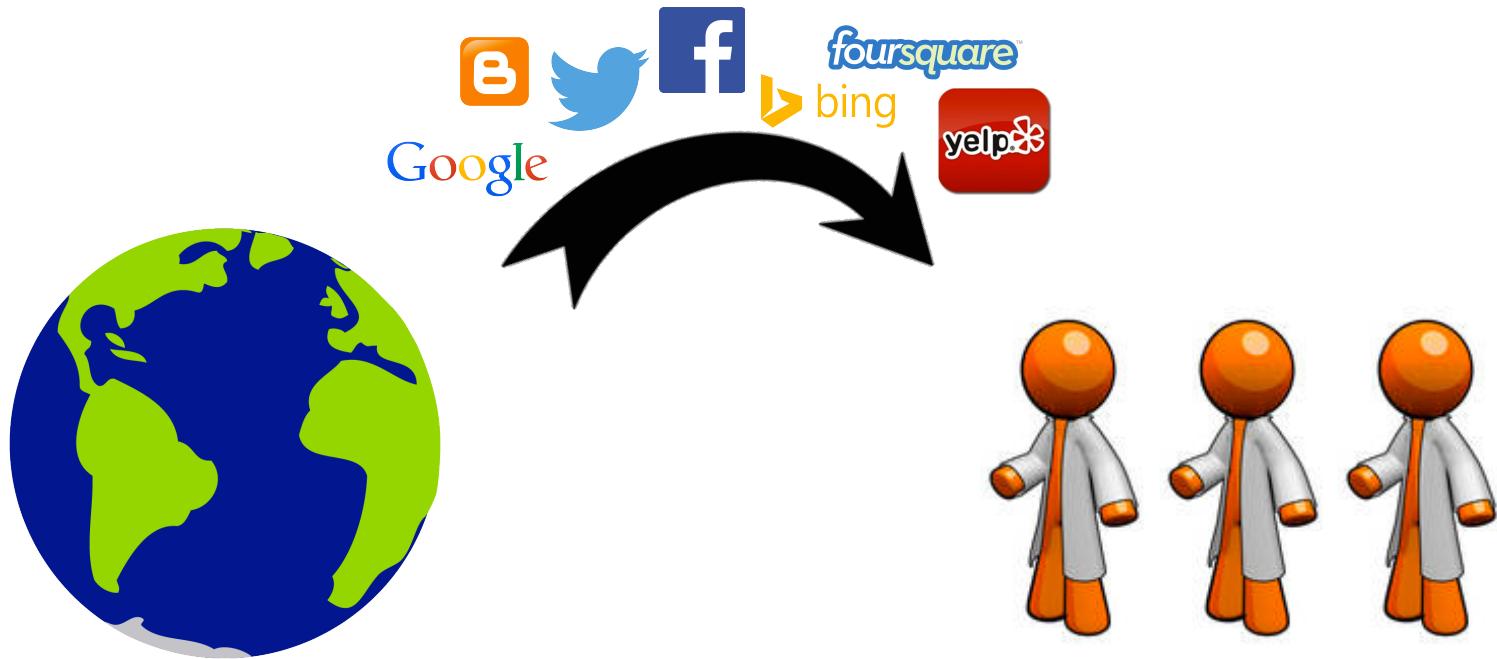
4.6K



5K



TEXT AS DATA



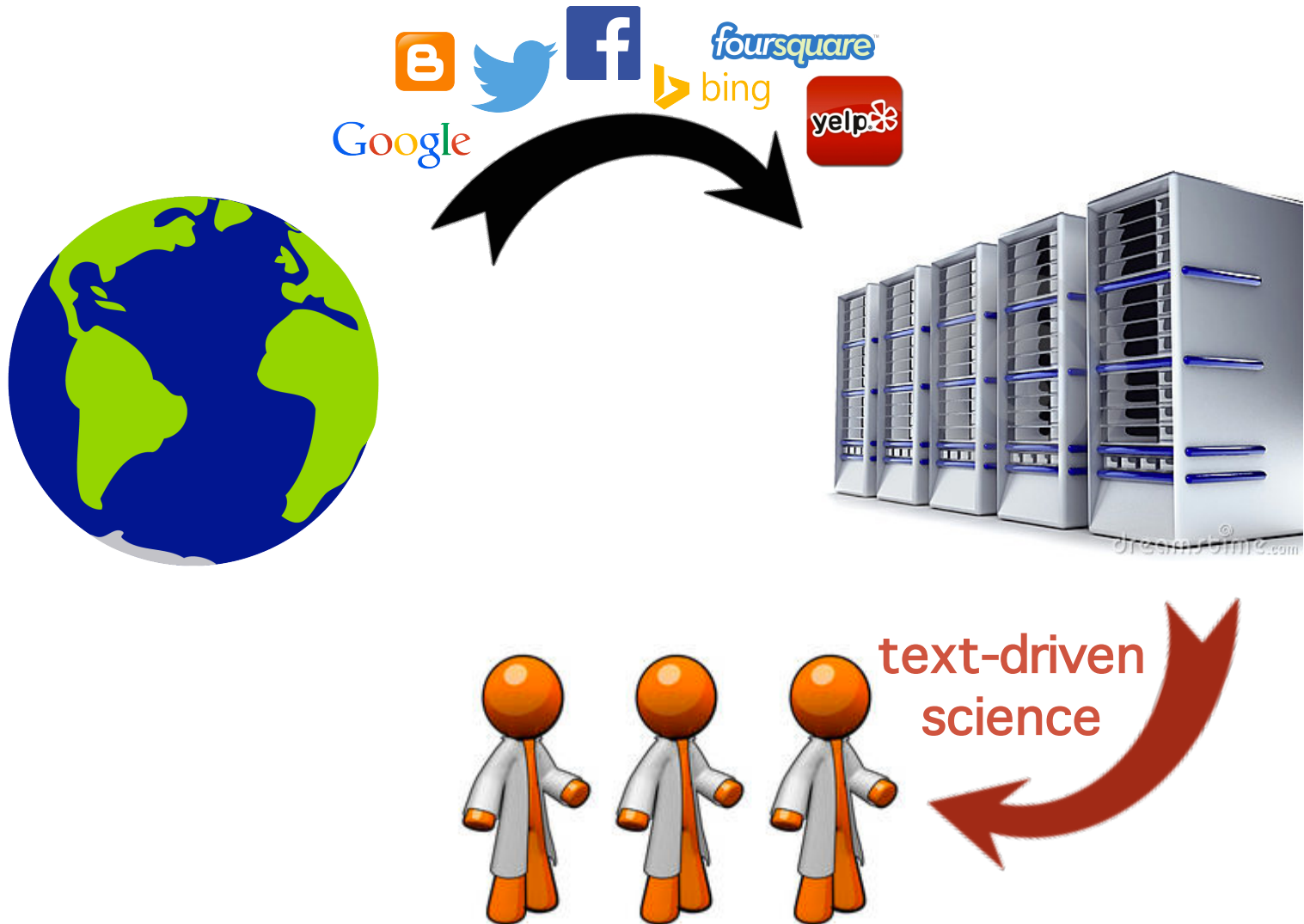
TEXT AS DATA



TEXT AS DATA



TEXT AS DATA



TEXT AS DATA



- Computational social science
- Computational journalism
- Crisis informatics
- Public health informatics
- Computational epidemiology

**text-driven
science**



TEXT AS DATA

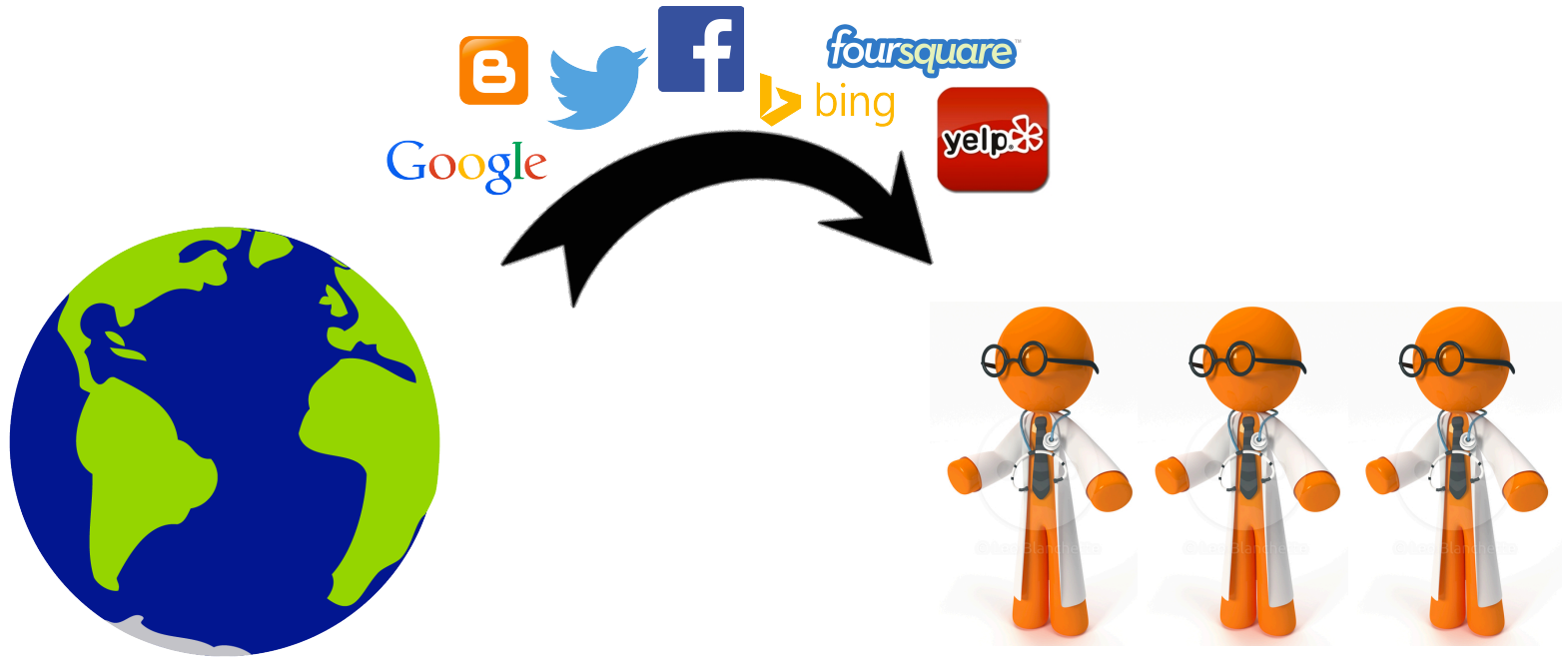


- Computational social science
- Computational journalism
- Crisis informatics
- Public health informatics
- Computational epidemiology

**text-driven
science**



TEXT-DRIVEN PUBLIC HEALTH

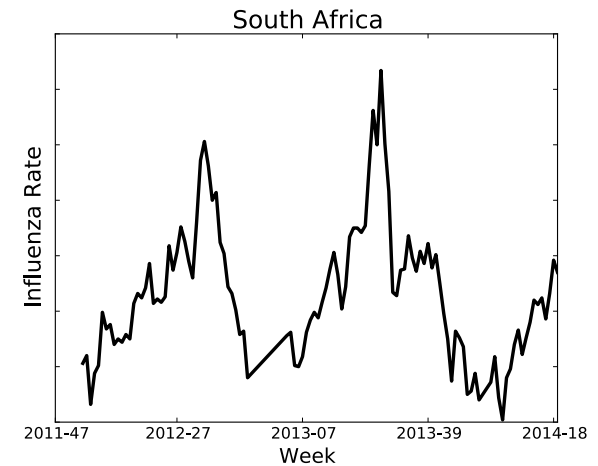
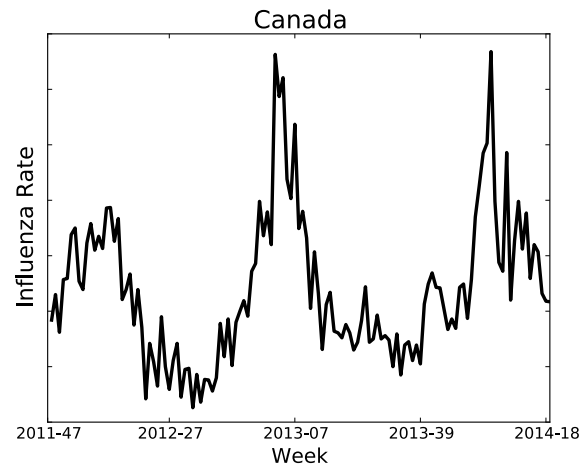
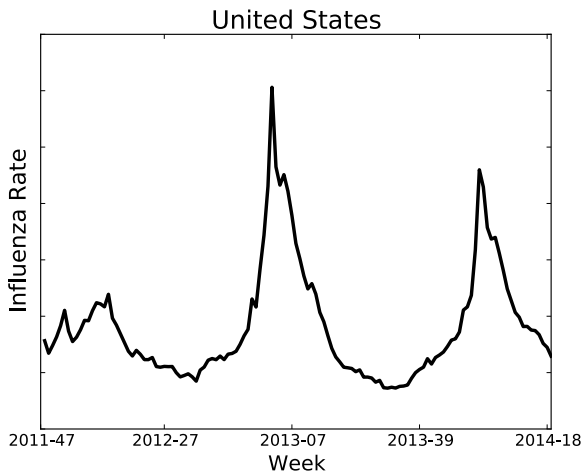


Let's look at some examples!

TEXT-DRIVEN PUBLIC HEALTH

FLU MONITORING

Important public health task: **disease surveillance**



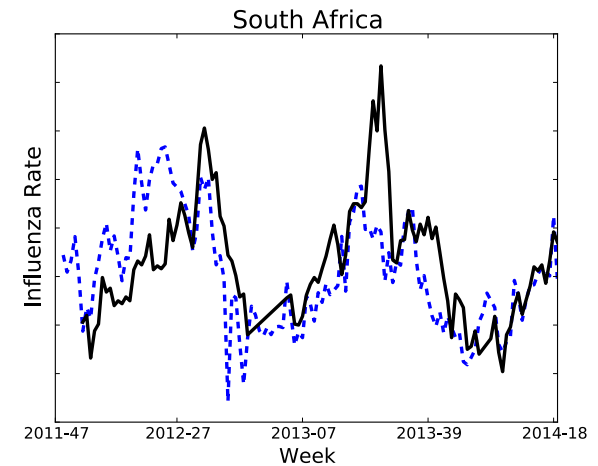
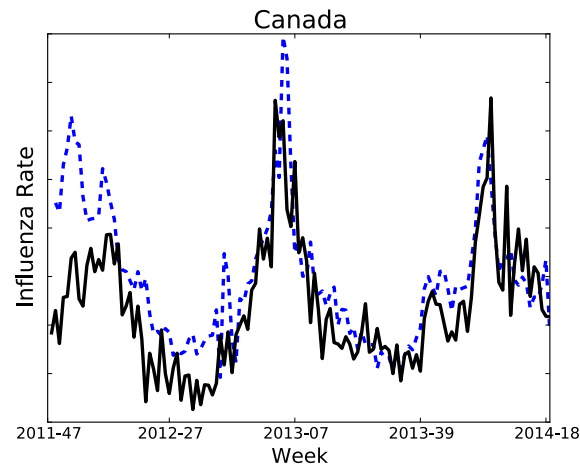
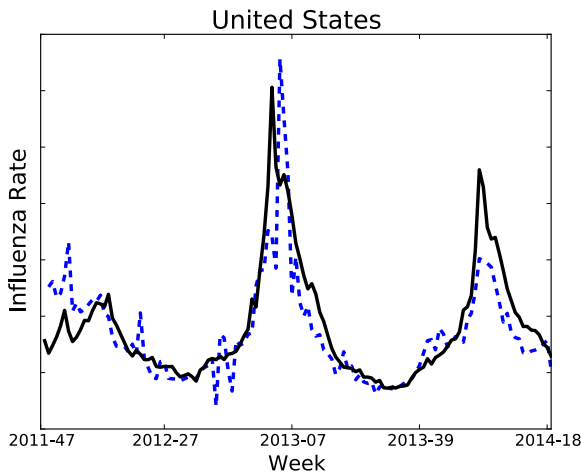
Many health agencies do flu monitoring

- But reports have a delay of ~2 weeks

TEXT-DRIVEN PUBLIC HEALTH

FLU MONITORING

Tracking the spread of influenza through tweets:



Paul, Dredze, Broniatowski (2014) **Twitter improves influenza forecasting.** *PLOS Currents: Outbreaks.*



Paul, Dredze, Broniatowski, Generous (2015) **Worldwide influenza surveillance through Twitter.** *AAAI Workshop on the World Wide Web and Public Health Intelligence.*



Lamb, Paul, Dredze (2013) **Separating fact from fear: Tracking flu infections on Twitter.** *NAACL.*



Broniatowski, Paul, Dredze (2013) **National and local influenza surveillance through Twitter: An analysis of the 2012-2013 influenza epidemic.** *PLOS ONE* 8(12): e83672.



TEXT-DRIVEN PUBLIC HEALTH

FLU MONITORING

Health Tweets

Home

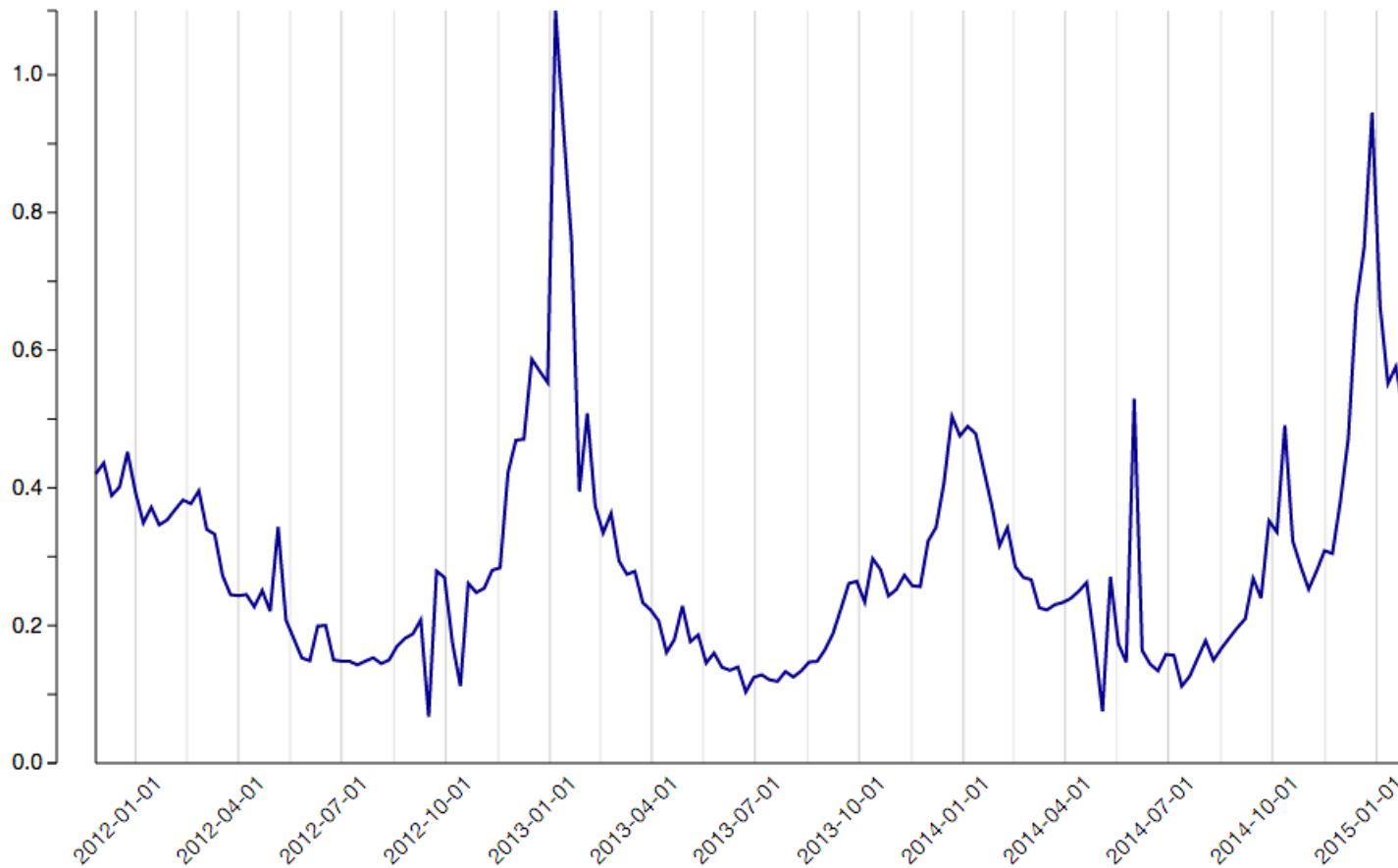
About

Trends

Flu Dashboard

My Account

Tracking Health Trends via Social Media



Legend

Influenza (United States)

[Modify plot](#)
[Modify Y Axis Range](#)
[Download to CSV](#)
[Create TinyUrl](#)

Air pollution in China

Public health tasks:

- Measure pollution levels
- Identify health effects
- Understand public response

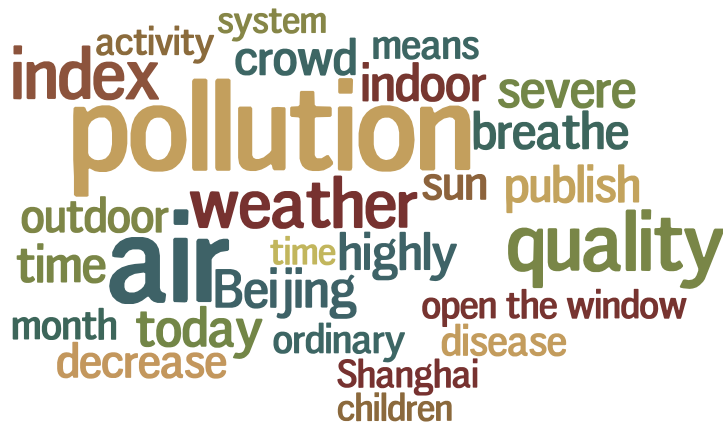
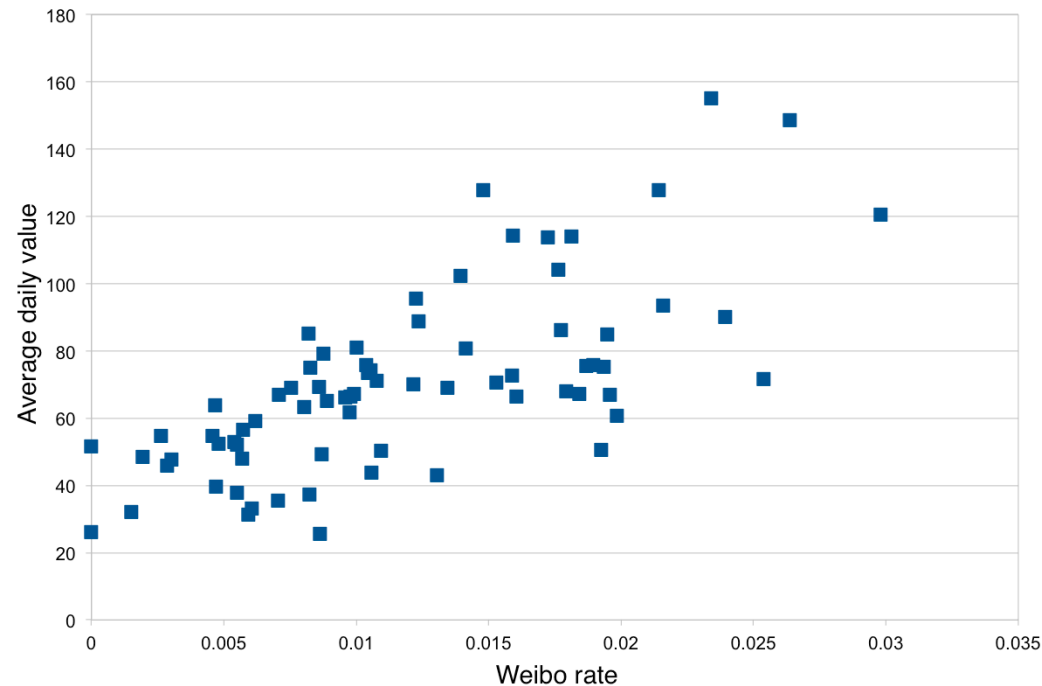


TEXT-DRIVEN PUBLIC HEALTH

AIR QUALITY

Monitoring air pollution through social media:

Relationship between pollution levels and weibos



Wang, Paul, Dredze (2015) Social media as a sensor of air quality and public response in China. *Journal of Medical Internet Research*.



TEXT-DRIVEN PUBLIC HEALTH

DRUG USE

New trends in **drug use**

- Record numbers of new drugs recently
- Health officials can be years behind



CNN News Video TV Opinions More...

U.S. World Politics Tech Health Entertainment Living Travel

Naked man chews off guy's face

May 29th, 2012
08:42 AM ET

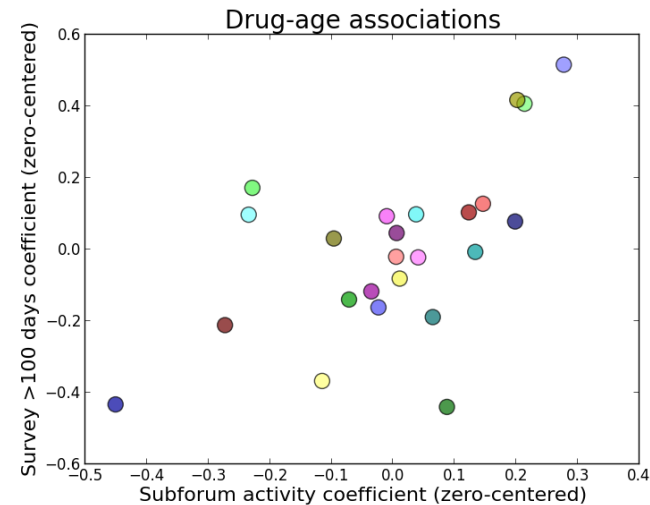
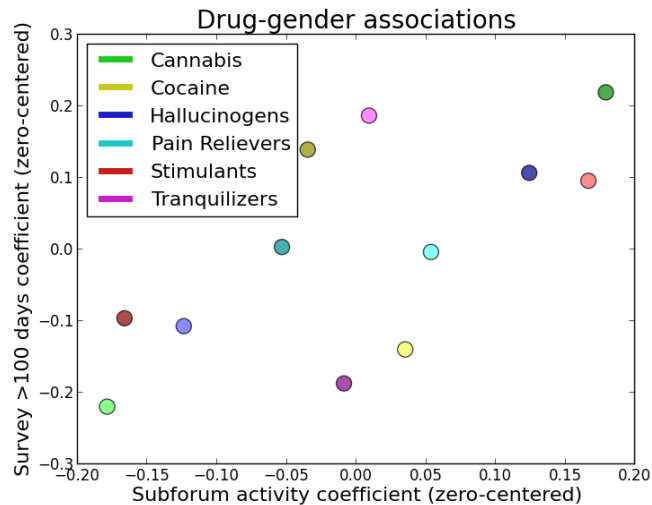
Reports: Miami 'zombie' attacker may have been using 'bath salts'


A naked man who chewed off the face of another man in what is being called a zombie-like attack may have been under the influence of "bath salts," a drug referred to as the new LSD, according to reports from CNN affiliates in Miami.

TEXT-DRIVEN PUBLIC HEALTH

DRUG USE

Analyzing online forums:



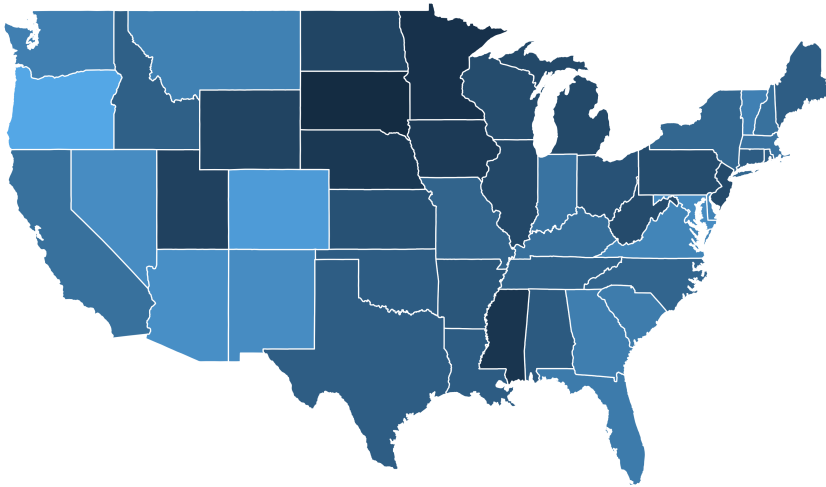
Paul, Dredze (2013) **Summarizing drug experiences with multi-dimensional topic models.** *North American ACL (NAACL).* 

Paul, Chisolm, Johnson, Vandrey, Dredze (in preparation) **Who participates in online drug communities? A large-scale analysis of demographic and temporal trends.** 

TEXT-DRIVEN PUBLIC HEALTH

HEALTHCARE QUALITY

Understanding **healthcare quality** from online reviews:



RateMDs.com Find a Doctor Bro

Latest Ratings:

"Wonderful bedside manner and extremely professional and honest! Highly recommend!"

[See this doctor's ratings](#)

Text from reviews is significantly predictive of external measures of healthcare quality

Paul, Wallace, Dredze (2013) **Analyzing online doctor ratings with a joint topic-sentiment model.** *AAAI Workshop on Expanding the Boundaries of Health Informatics Using AI.*



Wallace, Paul, Sarkar, Trikalinos, Dredze (2014) **A large-scale quantitative analysis of latent factors and sentiment in online doctor reviews.** *Journal of the American Medical Informatics Association* 21(6), 1098-1103.



TEXT-DRIVEN PUBLIC HEALTH

Many other applications:

- Air pollution in Chinese social media



Wang, Paul, Dredze (2015) **Social media as a sensor of air quality and public response in China.** *Journal of Medical Internet Research.*



- Health decision-making in search logs

Paul, White, Horvitz (2015) **Web search as medical decision support for cancer.** *WWW.*



- Public opinion in Twitter on public health issues:

- Gun control
- Vaccination
- Smoking

Benton, Paul, Hancock, Dredze (under review) **A joint model of topic and perspective in social media.**



TEXT-DRIVEN PUBLIC HEALTH

CBSNEWS Video US World Politics Entertainment Health MoneyWa

By MICHAEL CASEY / CBS NEWS / November 21, 2014, 6:00 AM

Scientists using social media to track air pollution in China

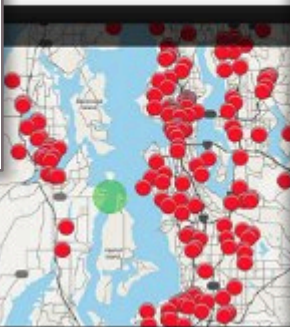


Clouds of smoke billow from a metal alloy factory in Gaolan county, Gansu province, northwest China AP



ice Entertainment Tech H

Mali Hillary Clinton Sarah P



SAVE BIG SUBSCRIBE TODAY



The Atlantic

POLITICS BUSINESS TECH ENTERTAINMENT HEALTH EDUCATION SEXES
JUST IN The Attention Machine PHOTO

10 Things We Can Learn From Your Health-Related Twitter Rants

HANS VILLARICA | JUL 15 2011, 12:46 PM ET



Unprecedented new research uses billions of tweets to reveal surprising patterns about cancer, obesity, and other ailments



Breaking news

South Korea makes new attempt to put satellite in orbit

As it faces hostility from neighboring North

Flu facts: Track 'em on Twitter

As the U.S. experiences its most severe flu seasons in years, res

The Washington Post

Search



Health & Science

Twitter becomes a tool for tracking flu epidemics and other public health issues



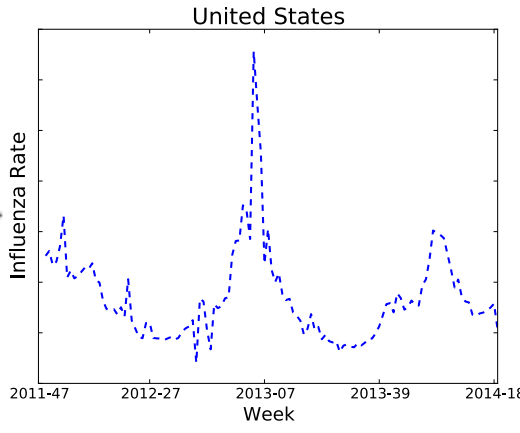
pressheretv.com viddler

TEXT AS DATA



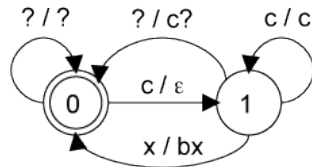
The screenshot shows a vertical scroll of tweets. Visible tweets include:

- Tiffany Thornton** (@heresTiffany) - 2h: Dear Lord, please let whatever Chris has be food poisoning n not the flu. My birthday is Saturday and I really don't want to spend it sick. (7 retweets, 43 likes)
- Roxeterawr** (@RoxeteraRbbons) - 2h: Oh no!! Everyone at work had flu and now I think I've got :!!!!!! 🙏 please noooooooooooooo I have too many things to do!! (15 retweets, 170 likes)
- Barack Obama** follows **Shahid Kamal Ahmad** (@shahidkamal) - 2h: 20 hour day, most of it work, one meal at 10:30pm, one toilet break in 13 hours and flu. And yet I'm ending the day totally psyched. (1 retweet, 4 likes)
- Richard Oliver** (@RichOliverActor) - 3h: Flu tabs taken & off to bed! leave you with another poem by Ricardo Pantelone as I head to my slumber #ActorLife x
if i were a tree,
what kind would i be?
a mammoth oak, tall and slender?
or a stumper version, silent, contemplative, tender?
a weeping willow i would be,
not for sadness, crying or misery,
for the willows roots lie strong and deep,
and by the winding rusby brook, in comfort sleep
Ricardo Pantelone (3 retweets, 9 likes)
- Jan Olsen** and 2 others follow **NFID** (@NFIDvaccines) - 3h: #FightFlu by adding prevention messages to lessons. Ready-to-use work plans available bit.ly/1us6yK7 #K12OS (3 retweets, 1 like)
- Kamran M. Riaz** (@kr156) - 3h: Can someone check if @bilimaher is down with the flu?As spokesperson for many atheists, his silence belies his approval #ChapelHillShooting (16 retweets, 12 likes)
- Kathleen Bachynski** and 10 others follow **CIDRAP** (@CIDRAP) - 4h: FLU SCAN: Parotitis in flu patients; Global flu update; H7N9 in China; Avian flu in Taiwan, Bulgaria ow.ly/1Up0Z (2 retweets, 1 like)
- NFID** (@NFIDvaccines) - 5h: Don't weather the #flu When flu hits, act fast! #FightFlu ow.ly/1M3u4 (3 retweets, 5 likes)
- Syndromic.org** (@ISDS) - 5h: HK's Dr. Ko Wing-man on Flu Reassortment Concerns (Avian Flu Diary) - bit.ly/1MASWc (3 retweets, 1 like)
- alex vespignani** and 11 others follow **Skoll Global** (@SkollGlobal) - 6h: Flu Near You featured on Fighting the Flu - FOX 8 WWUE New Orleans fox8live.com/Clip/1125348/... @FluNearYou (3 retweets, 1 like)
- NFID** (@NFIDvaccines) - 7h: Prompt use of antivirals is key this #flu season via @CDCFlu ow.ly/1Z2dQ (3 retweets, 1 like)
- Syndromic.org** (@ISDS) - 13h: US Flu Activity Down Slightly, but Elderly Hit Hard (CIDRAP) - bit.ly/1M83w3q (14 retweets, 5 likes)
- First** - 14h: Rabbits may be under swine flu, but claims to be div.it/BVYNbc (9 retweets, 3 likes)

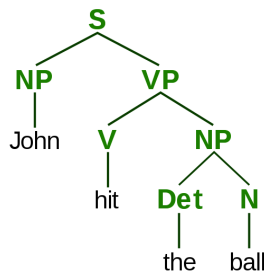


Structure of language:

- Morphology/strings



- Syntax/grammar



- Discourse/speech acts
- Topics/concepts

Paul, Eisner (2012) **Implicitly intersecting weighted automata using dual decomposition.** *NAACL*.

CS

Darling, Paul, Song (2012) **Unsupervised part-of-speech tagging in noisy domains with a syntactic-semantic Bayesian HMM.** *EACL Workshop on Social Media*.

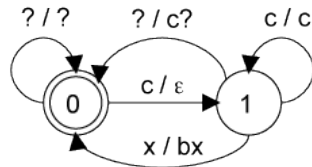
CS

Paul (2012) **Mixed membership Markov models for unsupervised conversation modeling.** *EMNLP*.

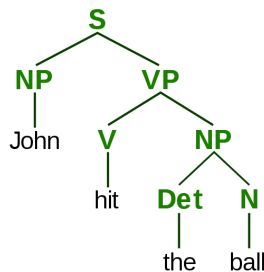
CS

Structure of language:

- Morphology/strings



- Syntax/grammar



- Discourse/speech acts
- Topics/concepts

Paul, Eisner (2012) **Implicitly intersecting weighted automata using dual decomposition.** *NAACL*.

CS

Darling, Paul, Song (2012) **Unsupervised part-of-speech tagging in noisy domains with a syntactic-semantic Bayesian HMM.** *EACL Workshop on Social Media*.

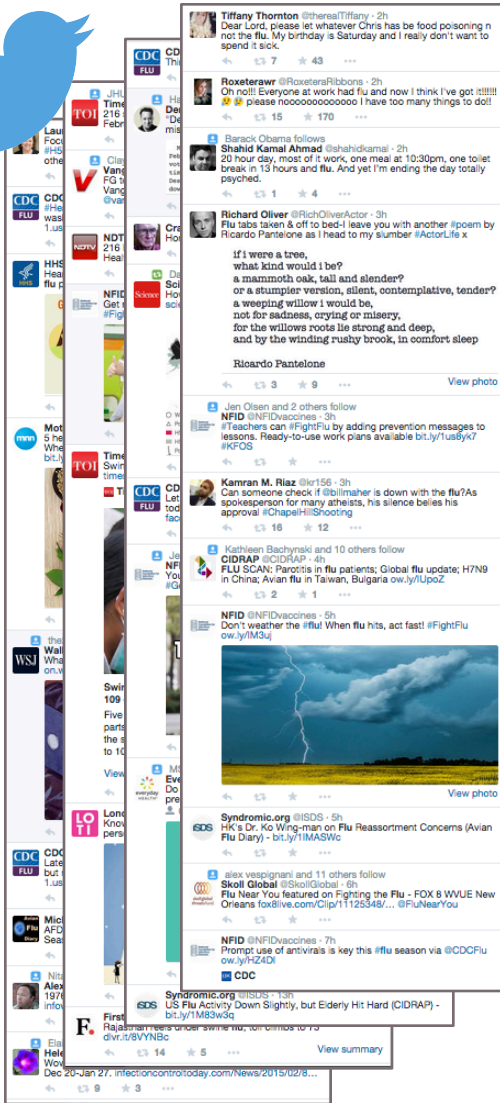
CS

Paul (2012) **Mixed membership Markov models for unsupervised conversation modeling.** *EMNLP*.

CS

TEXT AS DATA

NATURAL LANGUAGE PROCESSING



sick
sore
throat
feel
fever
flu

hurts
knee
ankle
hurt
neck
ouch


allergies
nose
eyes
allergy
allergic
sneezing

ow
teeth
tooth
wisdom
toothache
dentist

cancer
pray
surgery
hospital
pain
breast

body
pounds
gym
weight
lost
workout

Paul, Dredze (2011) You are what you tweet: Analyzing Twitter for public health. 5th International Conference on Weblogs and Social Media (ICWSM). 

Paul, Dredze (2014) Discovering health topics in social media using topic models. PLOS ONE 9(8). 

TOPIC MODELING

A topic model is a **statistical model** of text

- We pretend that our data (text) are the output of a probabilistic process that generates data

TOPIC MODELING

sick
sore
throat
feel
fever
flu

...

allergies
nose
eyes
allergy
allergic
sneezing

...

watch
watching
tv
killing
movie
seen

...

class
school
read
test
doing
finish

...

...

TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...

...

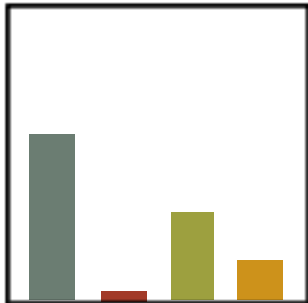
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



Michael Paul @mjp39 · Jan 24

I've had the flu and fever all week :(staying home from school and watching a lot of tv



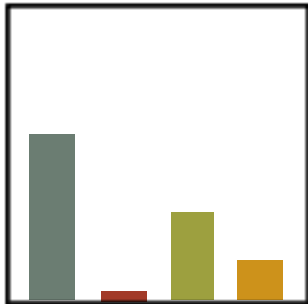
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



Michael Paul @mjp39 · Jan 24



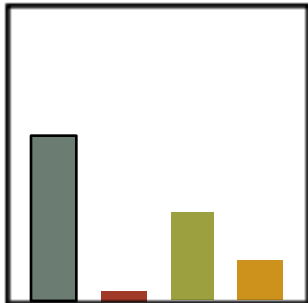
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



Michael Paul @mjp39 · Jan 24



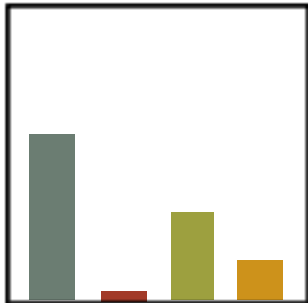
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



Michael Paul @mjp39 · Jan 24



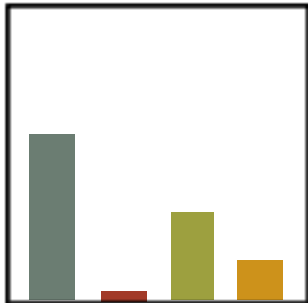
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



Michael Paul @mjp39 · Jan 24



fever



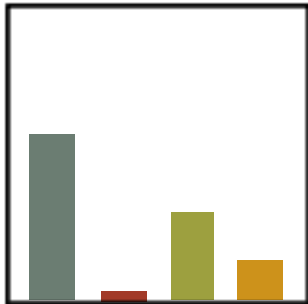
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



Michael Paul @mjp39 · Jan 24



fever



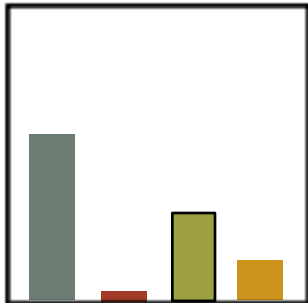
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



Michael Paul @mjp39 · Jan 24



fever



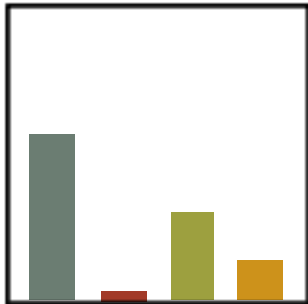
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



Michael Paul @mjp39 · Jan 24



fever



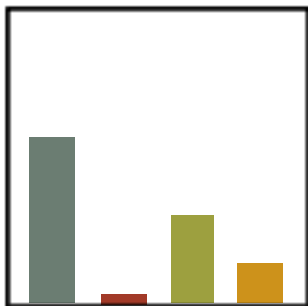
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



Michael Paul @mjp39 · Jan 24



fever

watching



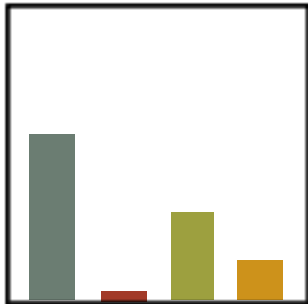
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



Michael Paul @mjp39 · Jan 24



fever

watching



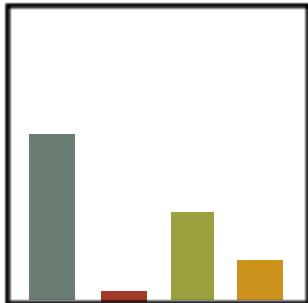
TOPIC MODELING

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



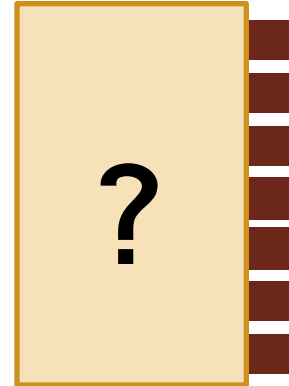
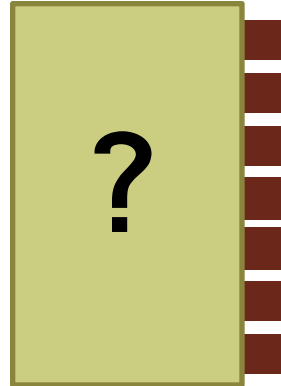
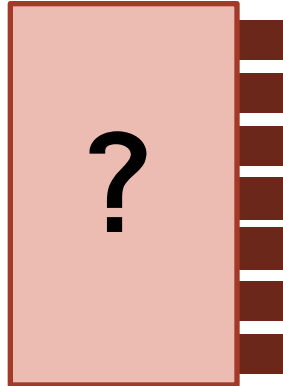
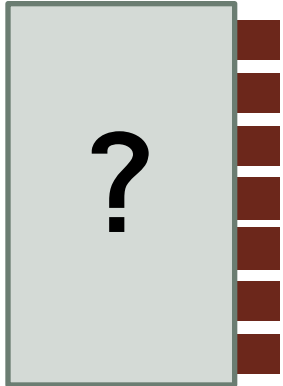
Michael Paul @mjp39 · Jan 24

I've had the flu and fever all week :(**staying** home from **school** and **watching** a lot of tv

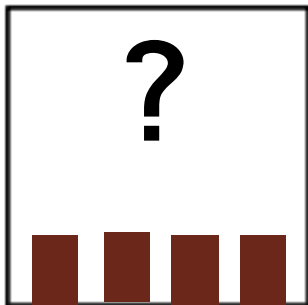




TOPIC MODELING

PARAMETER ESTIMATION







...



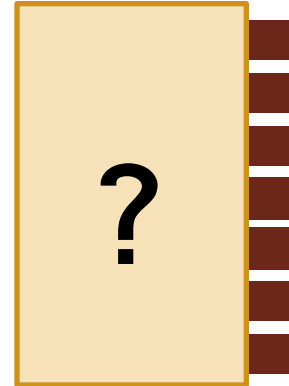
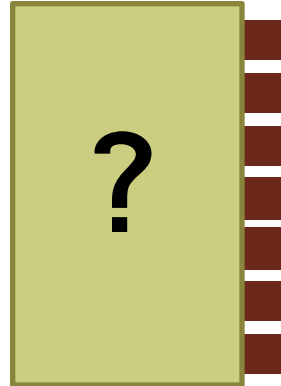
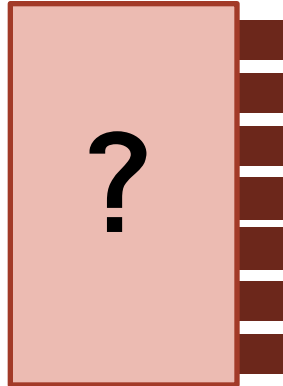
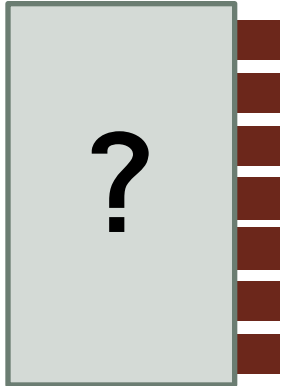
 **Michael Paul** @mjp39 · Jan 24 

I've had the flu and fever all week :(staying home from school and watching a lot of tv

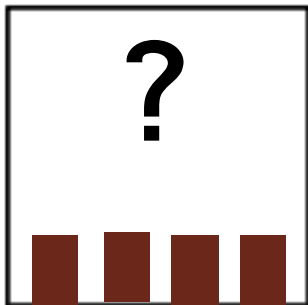
   

TOPIC MODELING

PARAMETER ESTIMATION





...



 **Michael Paul** @mjp39 · Jan 24

P (I've had the flu and fever all week :(staying home from school and watching a lot of tv) ?

TOPIC MODELING

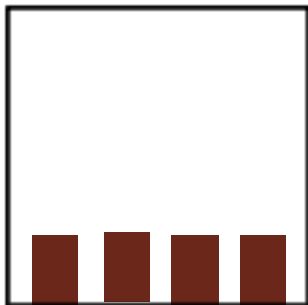
PARAMETER ESTIMATION


sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



P  **Michael Paul** @mjp39 · Jan 24
(I've had the flu and fever all week :(staying home from school and watching a lot of tv) ?



TOPIC MODELING

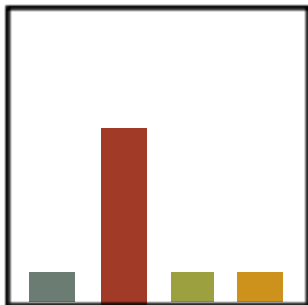
PARAMETER ESTIMATION


sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...



 **Michael Paul** @mjp39 · Jan 24
P (I've had the flu and fever all week :(staying home from school and watching a lot of tv) ?



TOPIC MODELING

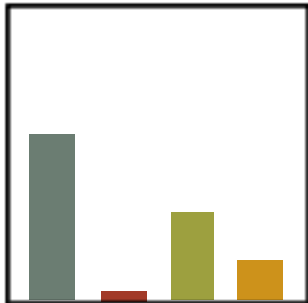
PARAMETER ESTIMATION


sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

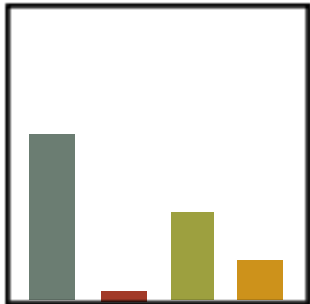
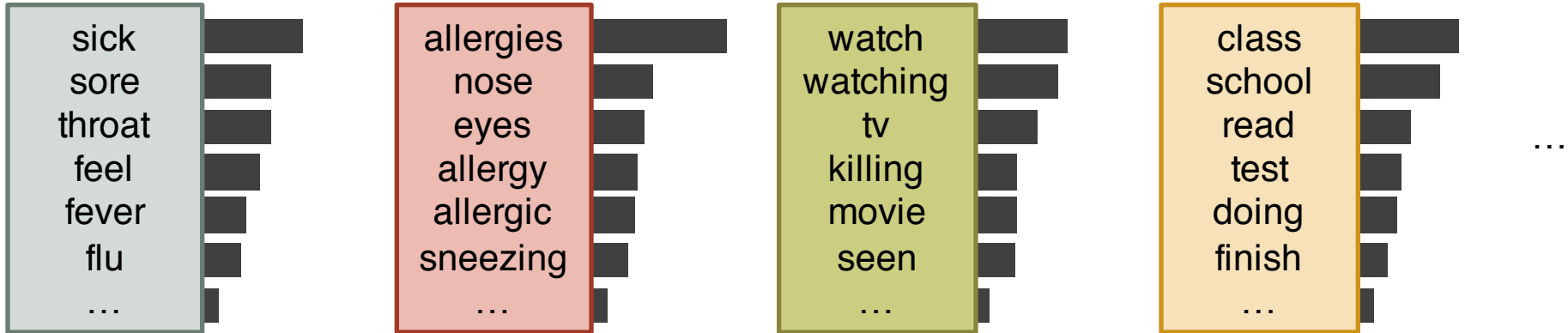
class
school
read
test
doing
finish
...



 **Michael Paul** @mjp39 · Jan 24
P (I've had the flu and fever all week :(staying home from school and watching a lot of tv) ?

TOPIC MODELING

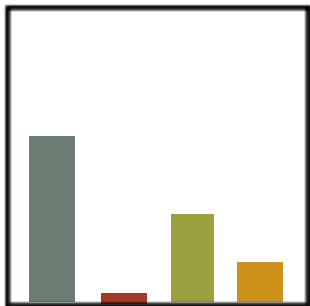
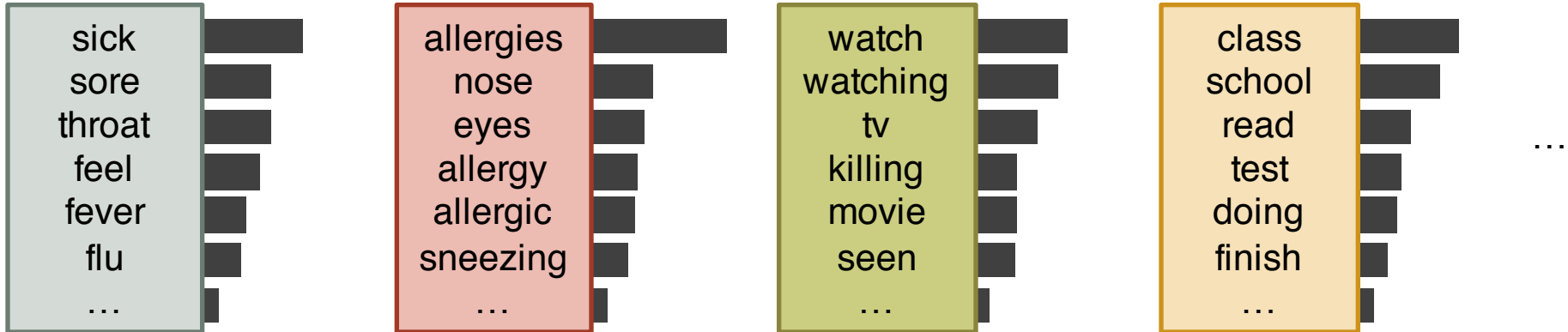
PRIORS



Our imaginary process also needs to generate all these distributions

TOPIC MODELING

PRIORS



Our imaginary process also needs to generate all these distributions

- We need a distribution over distributions
 - Called a **prior** distribution

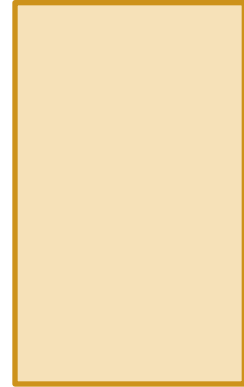
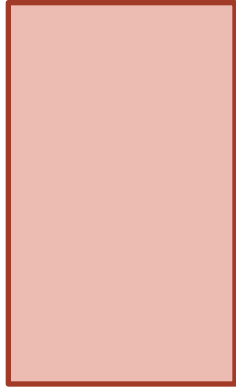
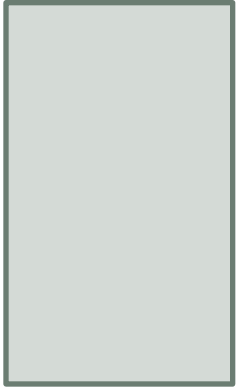
Dirichlet($\rho \times$ )

Precision

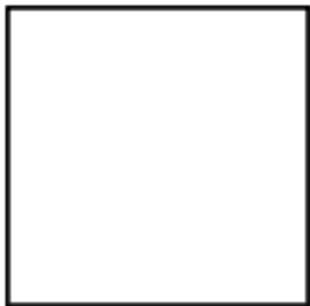
Mean

TOPIC MODELING

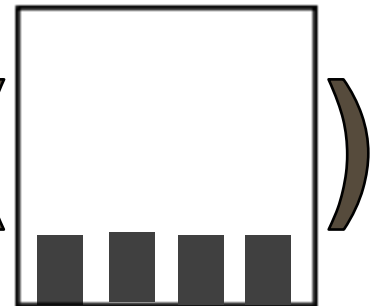
PRIORS



...



Dirichlet(



TOPIC MODELING

PRIORS

sick
sore
throat
feel
fever
flu
...

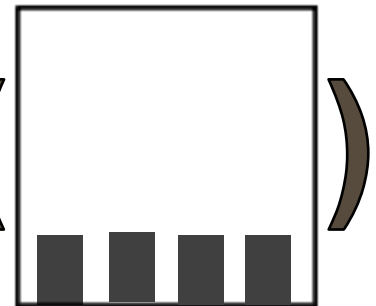
allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...

...

Dirichlet(



TOPIC MODELING

PRIORS

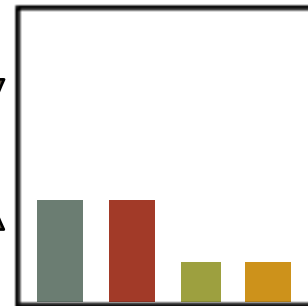
sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...

Dirichlet(



TOPIC MODELING

PRIORS

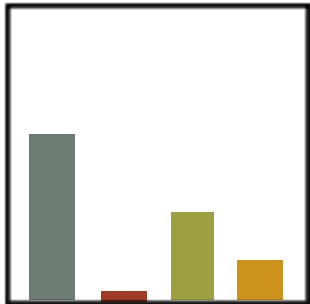
sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

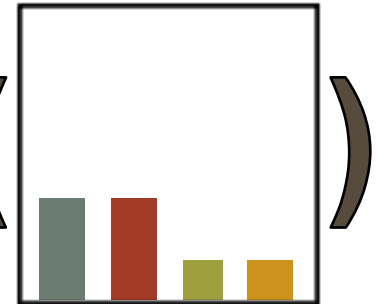
watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...

...

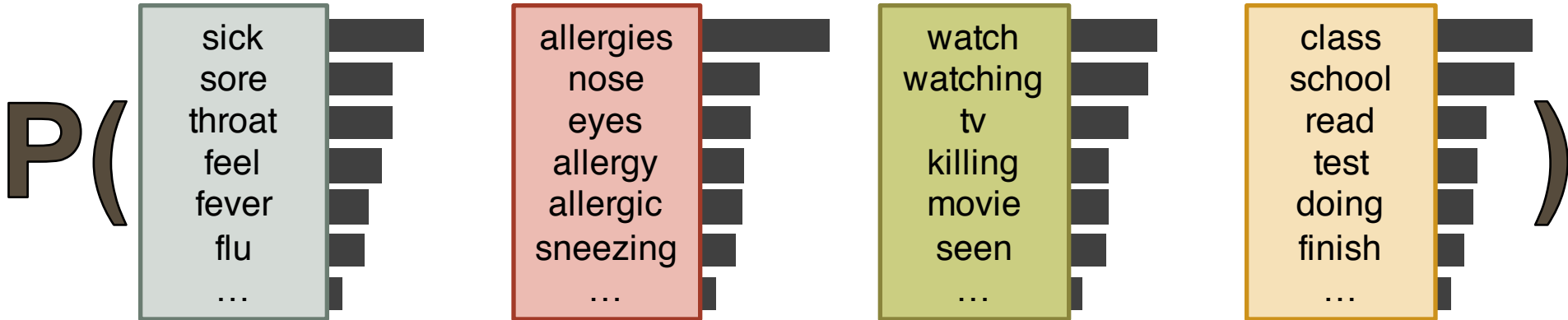


Dirichlet(

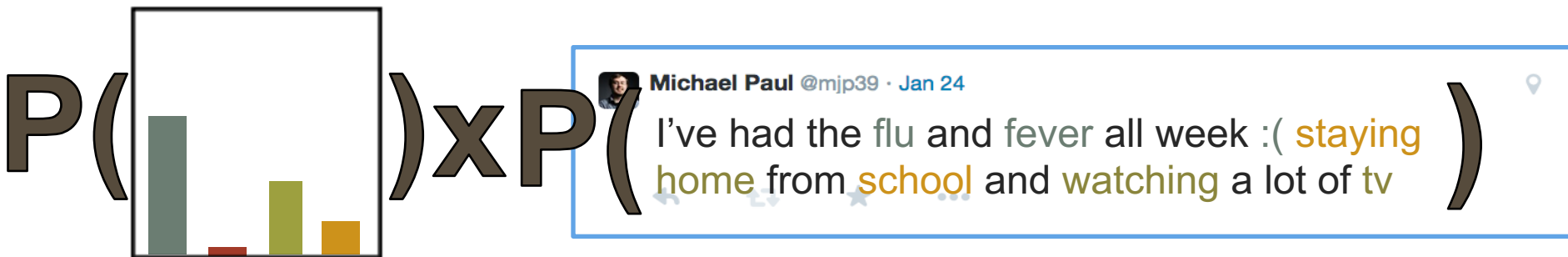


TOPIC MODELING

PRIORS



\times



Latent Dirichlet Allocation (LDA)

Blei, Ng, Jordan 2003

The topic and word distributions have Dirichlet priors

Latent Dirichlet Allocation (LDA)

Blei, Ng, Jordan 2003

The topic and word distributions have Dirichlet priors

Standard topic models are often insufficient for particular applications

- We need richer structure

TOPIC MODELING


APPLICATIONS

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...

Paul, Dredze (2011) **You are what you tweet: Analyzing Twitter for public health.** *5th International Conference on Weblogs and Social Media (ICWSM)*. 

Paul, Dredze (2014) **Discovering health topics in social media using topic models.** *PLOS ONE* 9(8): e103408. 

TOPIC MODELING

APPLICATIONS

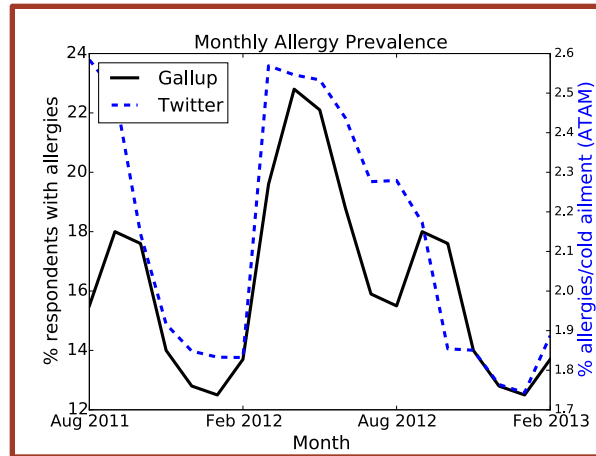
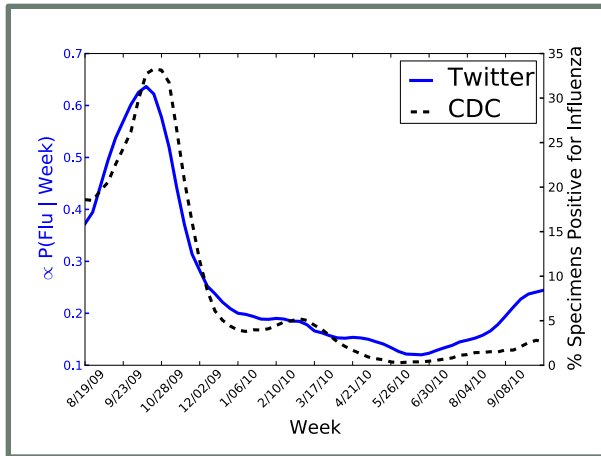
sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...

...



TOPIC MODELING

ADDING STRUCTURE

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...

About health issues

Irrelevant to health

...

TOPIC MODELING

ADDING STRUCTURE

sick
sore
throat
feel
fever
flu
...

allergies
nose
eyes
allergy
allergic
sneezing
...

watch
watching
tv
killing
movie
seen
...

class
school
read
test
doing
finish
...

About health issues

Irrelevant to health

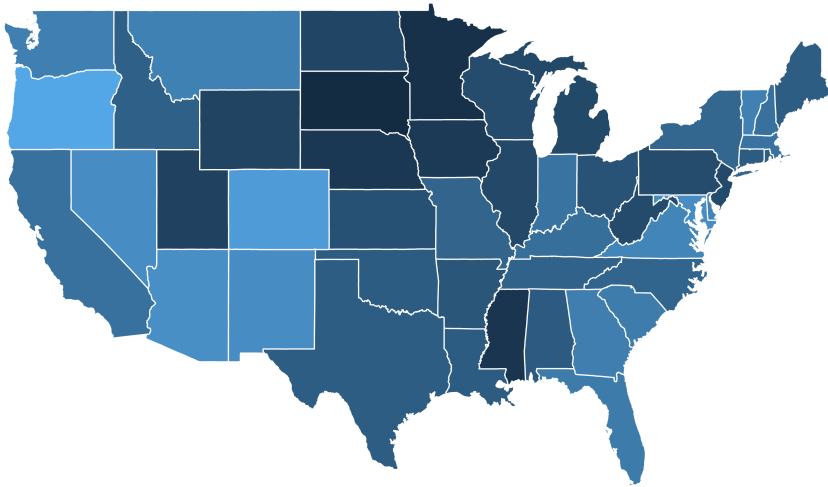
In general:

Topics can be organized in ways that are more interpretable to users

TOPIC MODELING

ADDING STRUCTURE

Understanding healthcare quality from online reviews:



RateMDs.com

Find a Doctor Bro

Latest Ratings:

"Wonderful bedside manner and extremely professional and honest! Highly recommend!"

[See this doctor's ratings](#)

TOPIC MODELING

ADDING STRUCTURE

Topics in online doctor reviews:



best
years
caring
care
patients
patient
recommend
family

time
staff
great
helpful
feel
questions
office
friendly

office
time
appointment
rude
staff
room
didn't
wait

TOPIC MODELING

ADDING STRUCTURE

Topics in online doctor reviews:



Both have **positive sentiment**

best
years
caring
care
patients
patient
recommend
family

time
staff
great
helpful
feel
questions
office
friendly

office
time
appointment
rude
staff
room
didn't
wait

Both about **staff/office issues**

TOPIC MODELING

ADDING STRUCTURE

Topics in online doctor reviews:



	Staff/Office	Personality	Surgery
Positive	time staff great helpful feel questions office friendly	best years caring care patients patient recommend family	surgery first son life surgeon daughter recommend thank
Negative	office time appointment rude staff room didn't wait	care medical patients doesn't help know don't problem	pain told went said surgery later didn't months

FACTORIAL LDA

A multi-dimensional topic model

Word distributions are grouped into different concepts

- e.g. **sentiment** and **aspect**

Paul and Dredze (2012) **Factorial LDA: Sparse multi-dimensional text models.**
Proceedings of *Advances in Neural Information Processing Systems (NIPS)*.

FACTORIAL LDA

DRUG DISCUSSIONS

Analyzing online drug forums:



3-dimensional model:

Drugs-Forum

- Drug type
- Route of administration (i.e. method of intake)
- Aspect

Drug (22 total)	Route	Aspect
<ul style="list-style-type: none">• Alcohol• Amphetamine• Cannabis• Cocaine• ...• Salvia• Tobacco	<ul style="list-style-type: none">• Injection• Oral• Smoking• Snorting	<ul style="list-style-type: none">• Chemistry• Culture• Effects• Health• Usage

Joint model with 3 factors:

- Drug type
- Route of administration (i.e. method of intake)
- Aspect

Each “topic” is a triple such as:

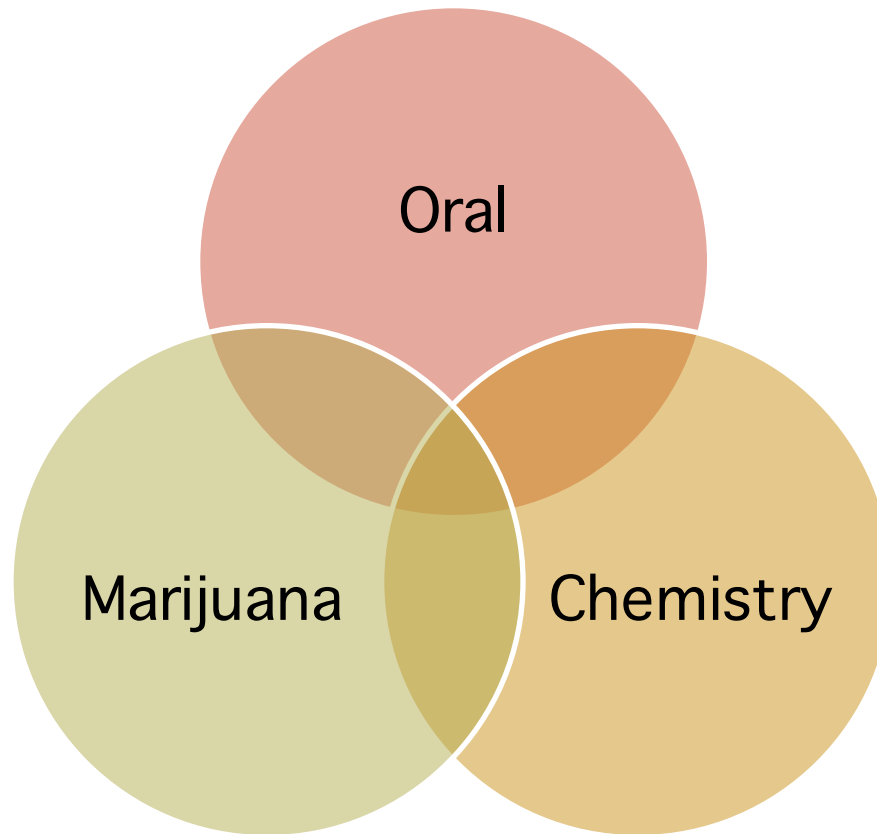
(Cocaine, Snorting, Health)

nose
pain
damage
blood
cocaine
problem

(Cocaine, Snorting, Usage)

coke
line
lines
nose
small
cut

Suppose we want to model: (Marijuana, Oral, Chemistry)



FACTORIAL LDA

DRUG DISCUSSIONS

Marijuana

weed
cannabis
thc
marijuana
stoned
bowl
bud
joint
blunt
herb
bong
pot
sativa
blaze
indica
smoking
blunts
...

Oral

capsules
consumes
toast
stomach
chewing
ambien
digestion
juice
absorbed
ingestion
meal
tiredness
chew
juices
gelatin
yogurt
fruit
...

Chemistry

solvent
extraction
evaporate
evaporated
solvents
evaporation
yield
chloride
alkaloids
tek
compounds
evaporating
atom
aromatic
non-polar
purified
jar
....

FACTORIAL LDA

DRUG DISCUSSIONS

Marijuana

weed
cannabis
thc
marijuana
stoned
bowl
bud
joint
blunt
herb
bong
pot
sativa
blaze
indica
smoking
blunts
...

Oral

capsules
consumes
toast
stomach
chewing
ambien
digestion
juice
absorbed
ingestion
meal
tiredness
chew
juices
gelatin
yogurt
fruit
...

Chemistry

solvent
extraction
evaporate
evaporated
solvents
evaporation
yield
chloride
alkaloids
tek
compounds
evaporating
atom
aromatic
non-polar
purified
jar
....

exp(

+

+

)

FACTORIAL LDA

DRUG DISCUSSIONS

Marijuana

weed
cannabis
thc
marijuana
stoned
bowl
bud
joint
blunt
herb
bong
pot
sativa
blaze
indica
smoking
blunts
...

Oral

capsules
consumes
toast
stomach
chewing
ambien
digestion
juice
absorbed
ingestion
meal
tiredness
chew
juices
gelatin
yogurt
fruit
...

Chemistry

solvent
extraction
evaporate
evaporated
solvents
evaporation
yield
chloride
alkaloids
tek
compounds
evaporating
atom
aromatic
non-polar
purified
jar
....

exp(

+

+

)

=

thc
method
extraction
plant
material
cannabis
simple
coffee
oil
contains
jar
dried
process
dry
water
extract
results
...

Dirichlet()

the
method
extraction
plant
material
cannabis
simple
coffee
oil
contains
jar
dried
process
dry
water
extract
results
...

FACTORIAL LDA

DRUG DISCUSSIONS

word distribution
for the triple:

(
 Marijuana
 Oral
 Chemistry
)

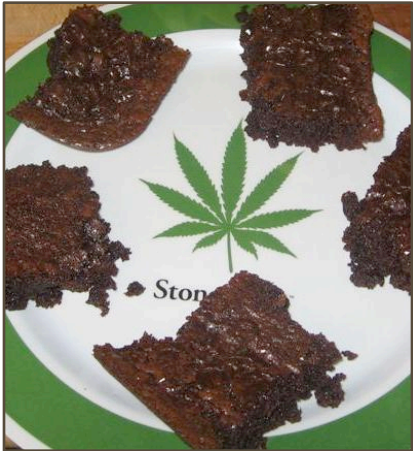
oil
water
butter
thc
weed
hash
cannabis
alcohol
make
milk
high
marijuana
add
...
mixture
hours
try
brownies

~ Dirichlet(

thc
the
method
extraction
plant
material
cannabis
simple
coffee
oil
contains
jar
dried
process
dry
water
extract
results
...

FACTORIAL LDA

DRUG DISCUSSIONS



oil
water
butter
thc
weed
hash
cannabis
alcohol
make
milk
high
marijuana
add
...
mixture
hours
try
brownies

~ Dirichlet()

thc
method
extraction
plant
material
cannabis
simple
coffee
oil
contains
jar
dried
process
dry
water
extract
results
...

FACTORIAL LDA

DRUG DISCUSSIONS

word distribution
for the triple:

(
 Marijuana
 Oral
 Chemistry
)

oil
water
butter
the
weed
hash
cannabis
alcohol
make
milk
high
marijuana
add
...
mixture
hours
try
brownies

~ Dirichlet(

the
method
extraction
plant
material
cannabis
simple
coffee
oil
contains
jar
dried
process
dry
water
extract
results
...

We can use this model to extract specific information about new drugs

- e.g. dosage, desired effects, negative effects *Drugs-Forum*

“What is the dosage when taking mephedrone orally?”

We can use this model to extract specific information about new drugs

- e.g. dosage, desired effects, negative effects *Drugs-Forum*

“What is the dosage when taking mephedrone orally?”

(
Mephedrone
Oral
Usage
)

We can use this model to extract specific information about new drugs

- e.g. dosage, desired effects, negative effects

Drugs-Forum

“What is the dosage when taking mephedrone orally?”

(
Mephedrone
Oral
Usage
)

If it is [someone who isn't you]'s first time using Mephedrone [someone who isn't me] recommends a 100mg oral dose on an empty stomach.

We can use this model to extract specific information about new drugs

- e.g. dosage, desired effects, negative effects

Drugs-Forum

“What is the dosage when taking mephedrone orally?”

(
Mephedrone
Oral
Usage
)

If it is [someone who isn't you]'s first time using Mephedrone [someone who isn't me] recommends a 100mg oral dose on an empty stomach.

Reference text:

It is recommended by users that Mephedrone be taken on an empty stomach. Doses usually vary between 100mg – 1g.

FACTORIAL LDA

SUMMARY

Word distributions are factored into multiple dimensions

Each topic prior is informed by parameters for each dimension

	Staff/Office	Personality	Surgery
Positive	time staff great helpful feel questions office friendly	best years caring care patients patient recommend family	surgery first son life surgeon daughter recommend thank
Negative	office time appointment rude staff room didn't wait	care medical patients doesn't help know don't problem	pain told went said surgery later didn't months

TOPIC STRUCTURES

FACTORIZATION

Components:

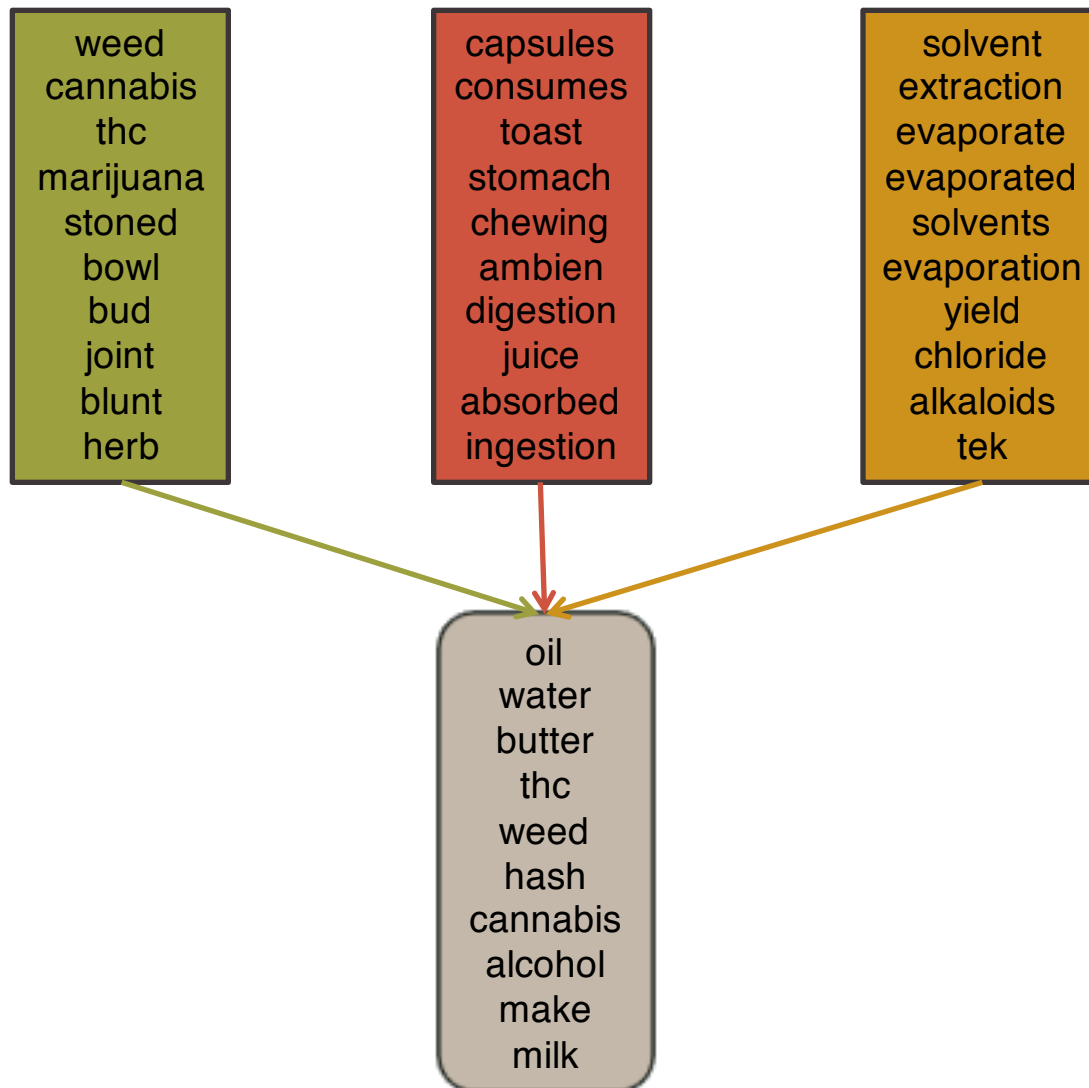
weed
cannabis
thc
marijuana
stoned
bowl
bud
joint
blunt
herb

capsules
consumes
toast
stomach
chewing
ambien
digestion
juice
absorbed
ingestion

solvent
extraction
evaporate
evaporated
solvents
evaporation
yield
chloride
alkaloids
tek

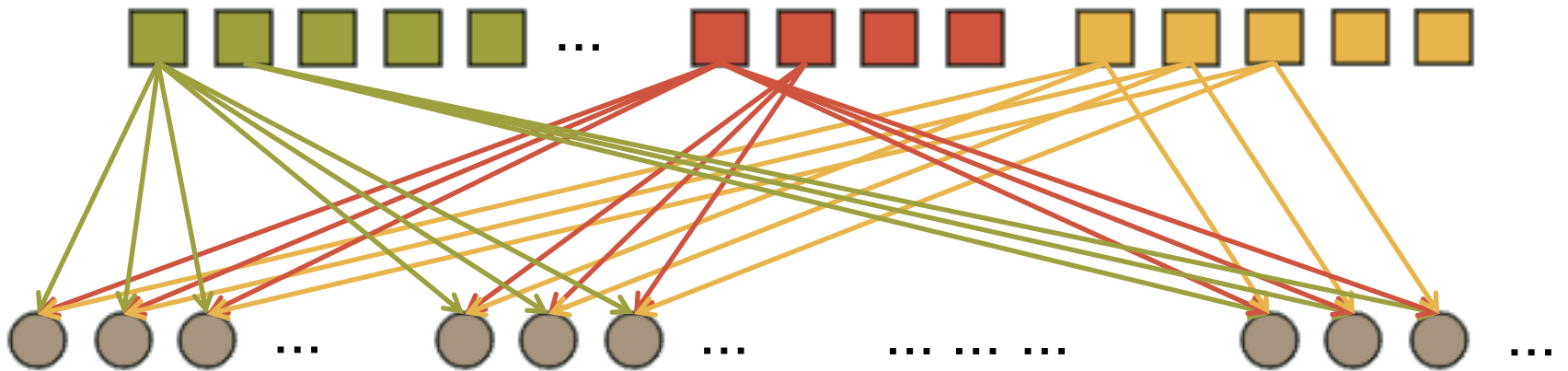
Topics:

oil
water
butter
thc
weed
hash
cannabis
alcohol
make
milk



TOPIC STRUCTURES

FACTORIZATION



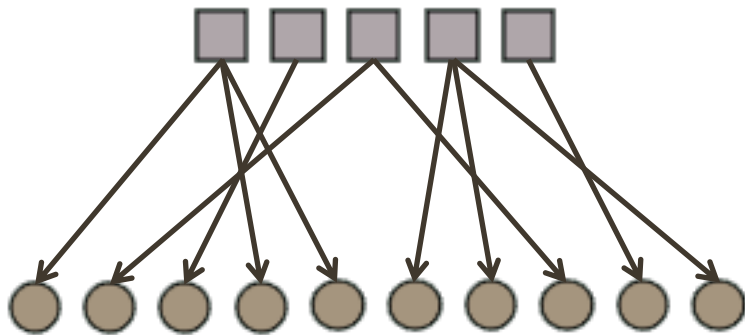
TOPIC STRUCTURES



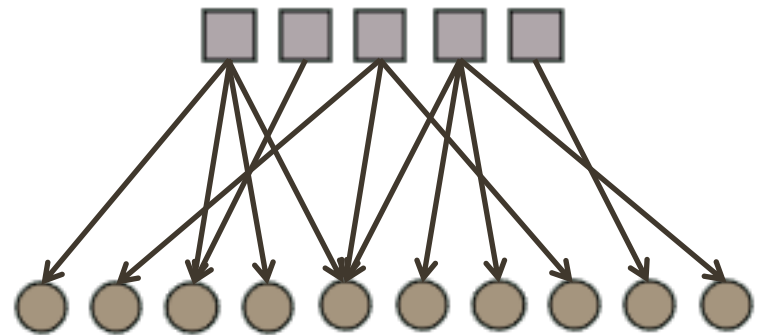
Generalization:

Components don't need to be factored into groups.
They can feed into topics in many different ways!

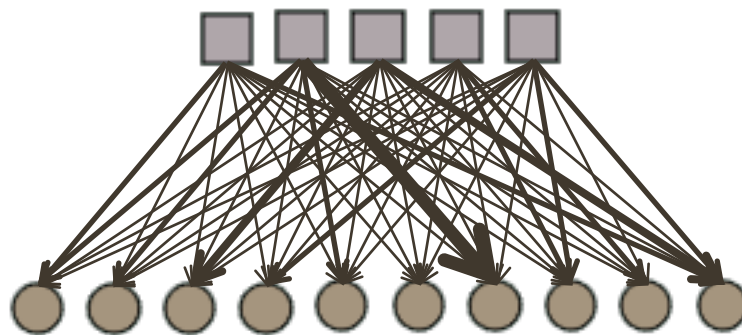
TOPIC STRUCTURES



Tree



Directed acyclic graph (DAG)



Weighted DAG

TOPIC STRUCTURES

TREE



“Surgery”

surgery
pain
went
dr
surgeon
told
procedure
months
performed
removed
left
fix
said
later
years

pain
surgery
dr
went
knee
foot
neck
mri
injury
shoulder
bone
months
told
surgeon
therapy

told
hospital
dr
blood
went
later
days
mother
said
er
cancer
weight
home
father
months

“Family”

dr
best
children
years
kids
cares
hes
care
old
daughter
child
husband
family
pediatrician
trust

dr
best
years
doctor
love
cares
ive
children
patients
hes
family
kids
seen
doctors
son

baby
son
pregnancy
dr
child
pregnant
ob
daughter
first
delivered
gyn
birth
delivery
section
hospital

TOPIC STRUCTURES

(SPARSE) DAG

“Surgery”

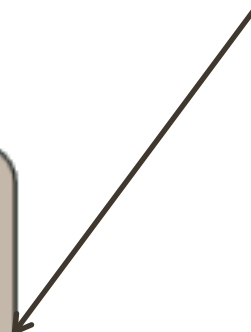
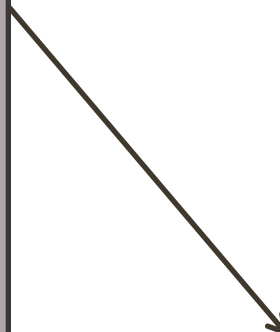
surgery
pain
went
dr
surgeon
told
procedure
months
performed
removed
left
fix
said
later
years



“Family”

dr
best
children
years
kids
cares
hes
care
old
daughter
child
husband
family
pediatrician
trust

dr
life
thank
saved
god
husband
heart
cancer
years
helped
doctors
hospital
father
man
able



SPRITE

Structured-prior topic models

A family of topic models in which the Dirichlet priors are functions of underlying components

Paul and Dredze (2015) **SPRITE: Generalizing topic models with structured priors.**
Transactions of the Association for Computational Linguistics (TACL) 3: 43-57.

SPRITE

SPRITE generalizes many existing topic models:

Model	Document priors	Topic priors
LDA	Single component	Single component
SCTM	Single component	Sparse binary β
SAGE	Single component	Sparse ω
FLDA	Binary δ is transpose of β	Factored binary β
PAM	α are supertopic weights	Single component
DMR	α are feature values	Single component

SPRITE

DEFINITION

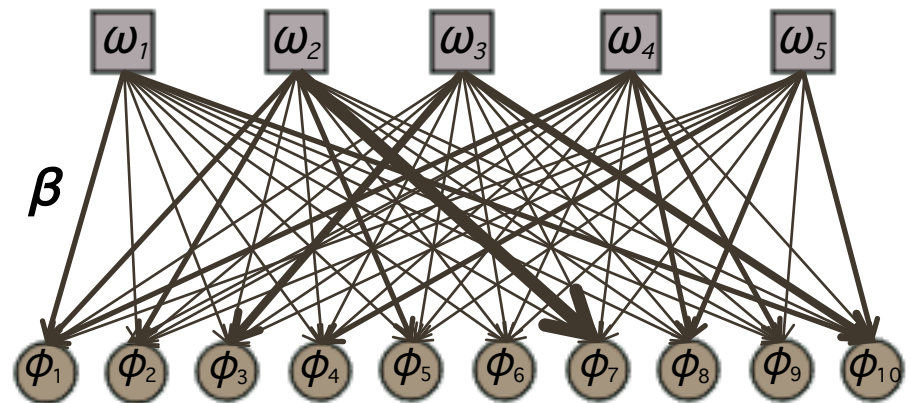
The priors over word distributions are weighted combinations of components:

$$\tilde{\phi}_{iv} = \exp\left(\sum_{c=1}^{C(\phi)} \beta_{ic} \omega_{cv}\right)$$

$$\phi_i \sim \text{Dirichlet}(\tilde{\phi}_i)$$



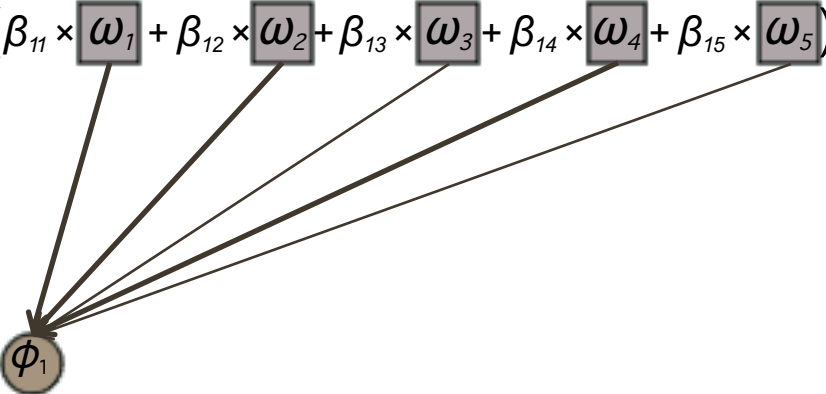
distribution over words in i th topic



SPRITE

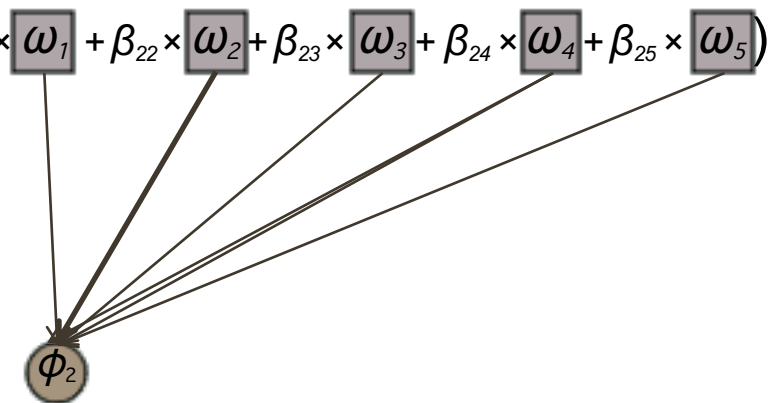
DEFINITION

The priors over word distributions are weighted combinations of components:

$$\exp(\beta_{11} \times \omega_1 + \beta_{12} \times \omega_2 + \beta_{13} \times \omega_3 + \beta_{14} \times \omega_4 + \beta_{15} \times \omega_5)$$


The diagram illustrates the relationship between the parameters in the equation and a central node. A circular node labeled ϕ_1 is positioned below the equation. Five arrows originate from this node and point to the terms ω_1 , ω_2 , ω_3 , ω_4 , and ω_5 within the equation, indicating that these terms are components of the prior distribution.

The priors over word distributions are weighted combinations of components:

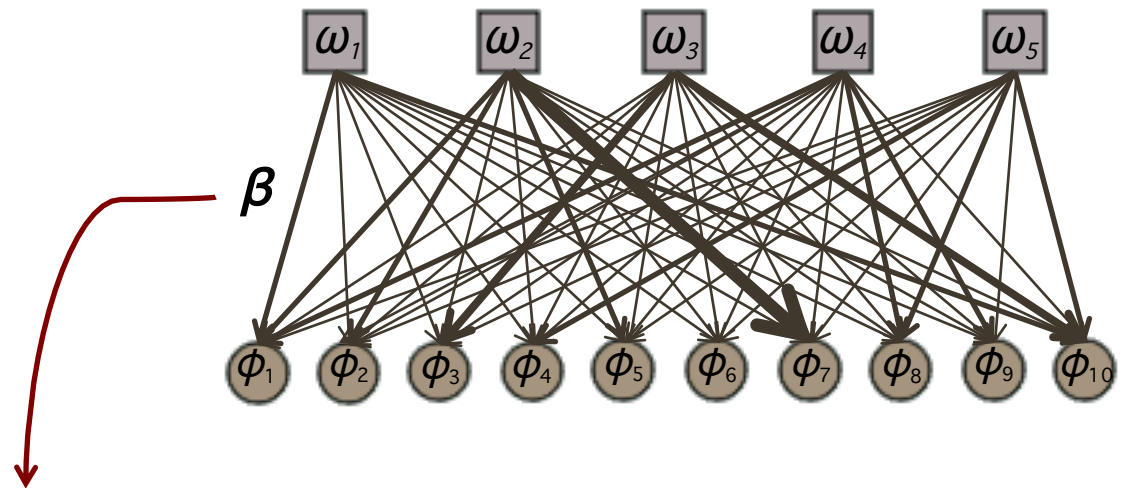
$$\exp(\beta_{21} \times \omega_1 + \beta_{22} \times \omega_2 + \beta_{23} \times \omega_3 + \beta_{24} \times \omega_4 + \beta_{25} \times \omega_5)$$


The diagram illustrates the mathematical expression above. A central circular node labeled ϕ_2 is connected by lines to five square nodes labeled ω_1 , ω_2 , ω_3 , ω_4 , and ω_5 . These nodes are arranged in a horizontal line above ϕ_2 . The lines represent the weights β_{21} through β_{25} in the equation.

SPRITE

DEFINITION

The priors over word distributions are weighted combinations of components:

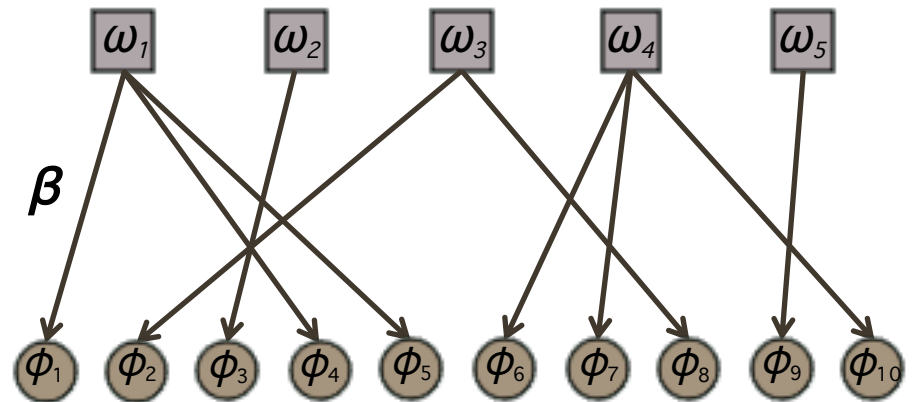


We can induce different structures by constraining the values of β

SPRITE

DEFINITION


The priors over word distributions are weighted combinations of components:



SPRITE

DEFINITION

The priors over word distributions are weighted combinations of components:

$$\exp(1 \times \omega_1 + 0 \times \omega_2 + 0 \times \omega_3 + 0 \times \omega_4 + 0 \times \omega_5)$$


A diagram illustrating the relationship between the weight ω_1 and the component ϕ_1 . An arrow points from the ω_1 term in the equation above to a circular node labeled ϕ_1 .

SPRITE

DEFINITION

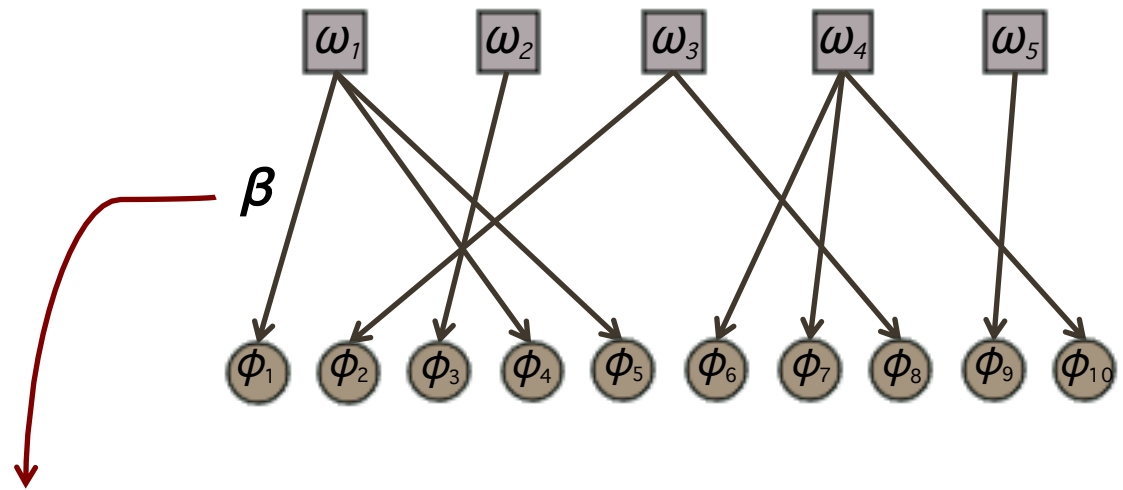
The priors over word distributions are weighted combinations of components:

$$\exp(0 \times \omega_1 + 0 \times \omega_2 + 1 \times \omega_3 + 0 \times \omega_4 + 0 \times \omega_5)$$



ϕ_2

The priors over word distributions are weighted combinations of components:

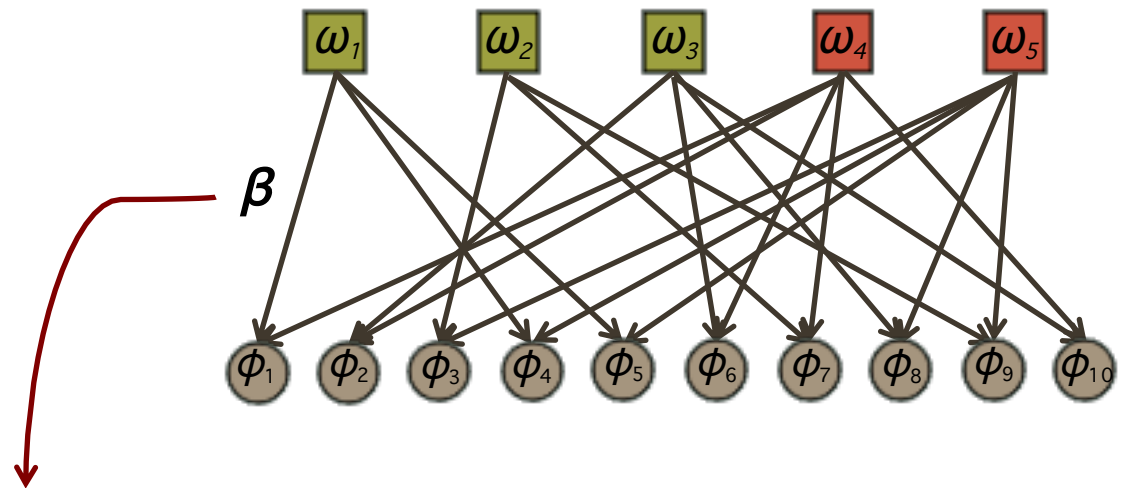


Tree: Each topic's β vector is zero in all but one component

SPRITE

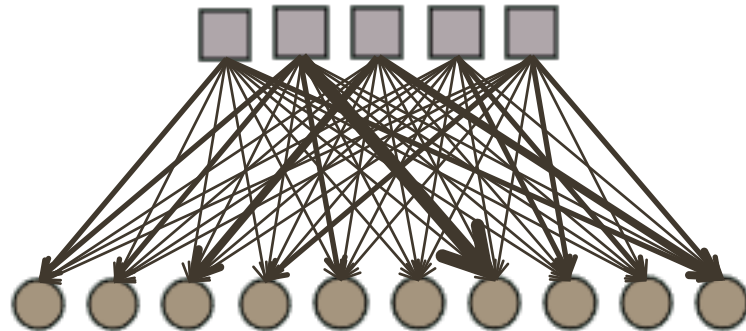
DEFINITION

The priors over word distributions are weighted combinations of components:



Factorization: Like a tree, but a nonzero component in each factor

- Organizes topics in a variety of useful ways
 - Can be tailored toward different applications
- Generalizes many topic models
 - While opening up new possibilities
- Allows practitioners to make sense of big text data
 - Can drive new scientific research



MOVING FORWARD



MOVING FORWARD



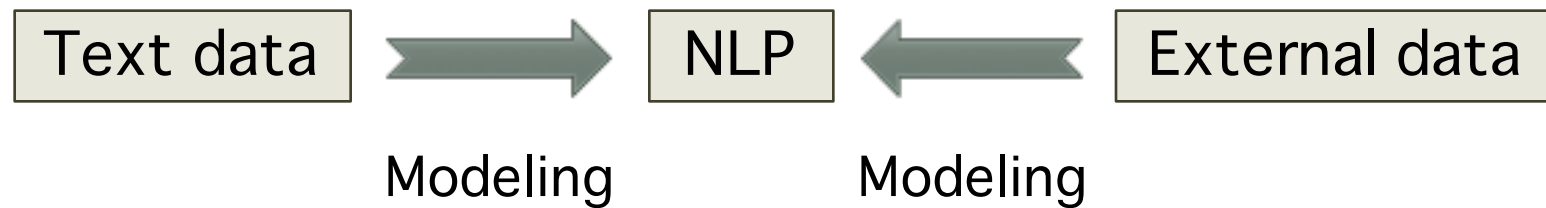
MOVING FORWARD

BETTER MODELS



MOVING FORWARD

BETTER MODELS



MOVING FORWARD

BETTER MODELS



GET STARTED

SEARCH OVER [138,198 DATASETS](#)



Agriculture



Business



Climate



Consumer



Ecosystems



Education



Energy



Finance



Health



Local
Government



Manufacturing



Ocean



Public Safety

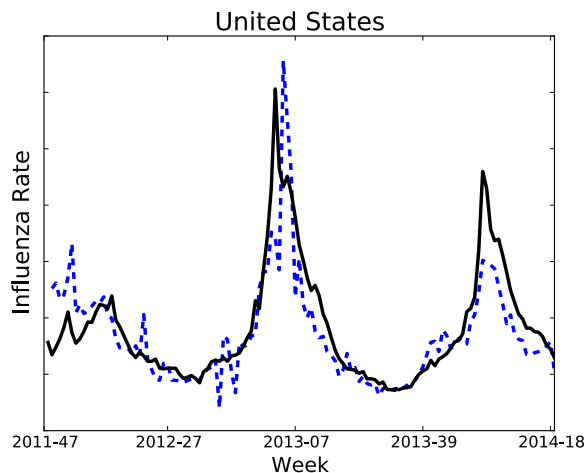


Science &
Research

Challenges with big models:

- Spurious correlations between **text** and **datasets**

Modeling flu prevalence in Twitter:

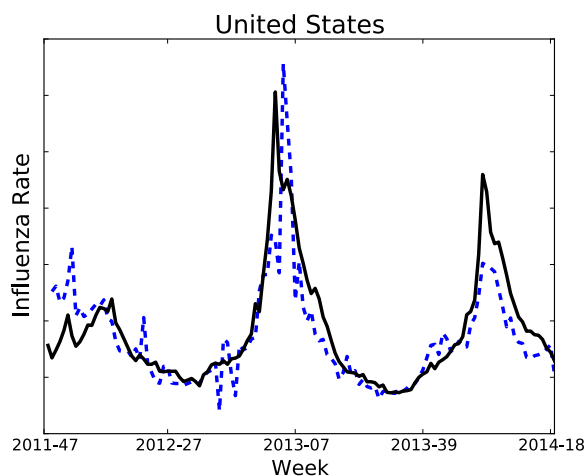


flu
sick
swine
shot
cancer
fever
h1n1
#beatcancer
better
getting
home
halloween
breast

Challenges with big models:

- Spurious correlations between **text** and **datasets**

Modeling flu prevalence in Twitter:



flu
sick
swine
shot
cancer
fever
h1n1
#beatcancer
better
getting
home
halloween
breast

Covariance
prior:



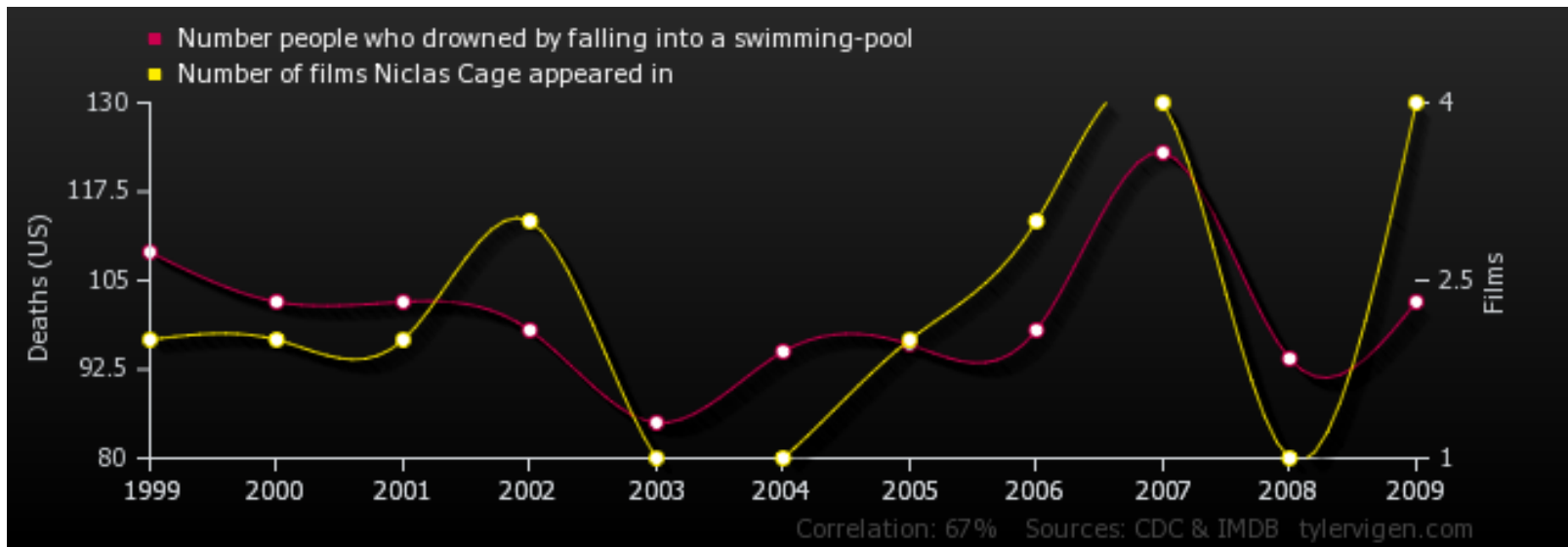
sick
flu
swine
better
shot
getting
cancer
home
hope
fever
feel
feeling
h1n1

MOVING FORWARD

BETTER MODELS

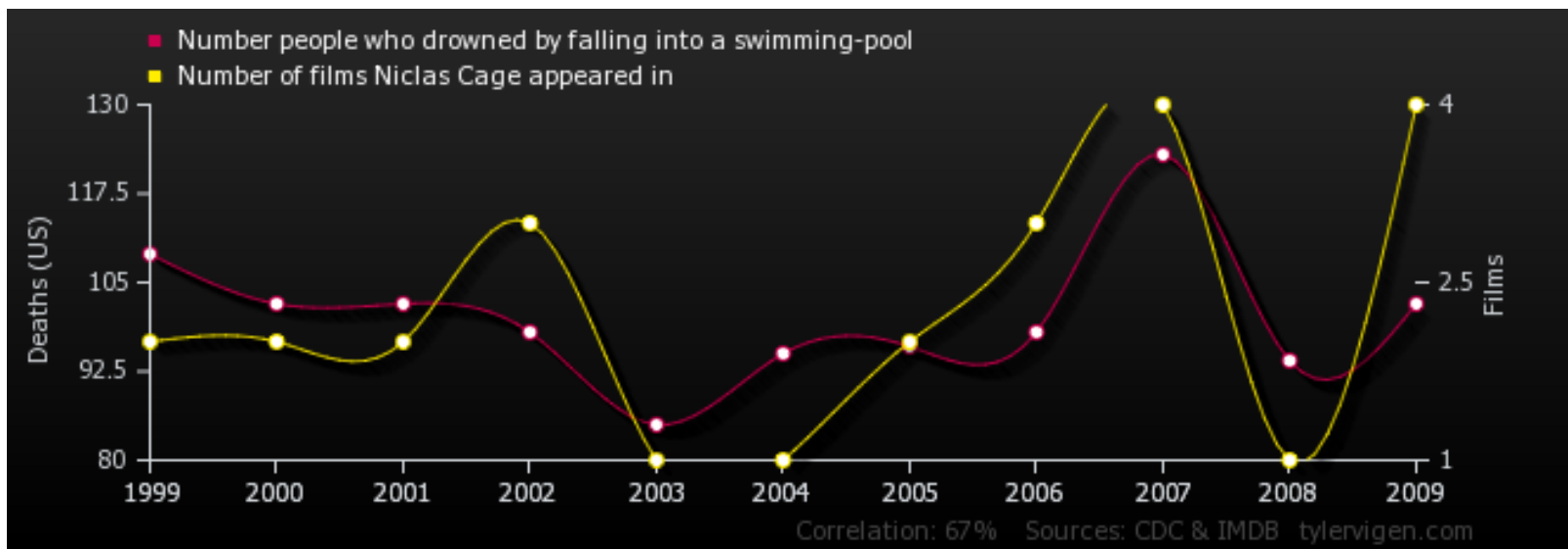
Challenges with big models:

- Spurious correlations between different datasets



Challenges with big models:

- Spurious correlations between different datasets



- Need models with domain expertise
 - Opportunities for interactive machine learning

Challenges with big models:

- **Solutions:** language structure and human feedback

MOVING FORWARD



MOVING FORWARD

NEW QUESTIONS

What happens when this article gets published?



MOVING FORWARD

NEW QUESTIONS

What happens when this article gets published?



Most cancers are due to bad luck, not lifestyle, researchers say.
time.com/3651785/cancer... ...Anyone fancy a smoke?

Most cancer just 'bad luck'?
time.com/3651785/cancer... So I'm going for checkup w/ my magic eightball today.
Just saved copay

MOVING FORWARD

NEW QUESTIONS

Variation in cancer risk among tissues can be explained by the number of stem cell divisions

Cristian Tomasetti^{1*} and Bert Vogelstein^{2*}

Some tissue types give rise to human cancers millions of times more often than other tissue types. Although this has been recognized for more than a century, it has never been explained. Here, we show that the lifetime risk of cancers of many different types is strongly correlated (0.81) with the total number of divisions of the normal self-renewing cells maintaining that tissue's homeostasis. These results suggest that only a third of the variation in cancer risk among tissues is attributable to environmental factors or inherited predispositions. The majority is due to "bad luck," that is, random mutations arising during DNA replication in normal, noncancerous stem cells. This is important not only for understanding the disease but also for designing strategies to limit the mortality it causes.

Extrême variation in cancer incidence across different tissues is well known; for example, the lifetime risk of being diagnosed with cancer is 6.9% for lung, 1.08% for thyroid, 0.6% for brain and the rest of the nervous system, 0.003% for pelvic bone and 0.00072% for laryngeal cartilage (1-5). Some of these differences are associated with well-known risk factors such as smoking, alcohol use, ultraviolet light, or human papilloma virus (HPV) (6, 5), but this applies only to specific populations

exposed to potent mutagens or viruses. And such exposures cannot explain why cancer risk in tissues within the alimentary tract can differ by as much as a factor of 24 (esophagus (0.51%), large intestine (4.82%), small intestine (0.20%), and stomach (0.86%)) (3). Moreover, cancers of the small intestinal epithelium are three times less common than brain tumors (3), even though small intestinal epithelial cells are exposed to much higher levels of environmental mutagens than are cells within the brain, which are protected by the blood-brain barrier.

Another well-studied contributor to cancer is



Bad Luck of Random Mutations Plays Predominant Role in Cancer, Study Shows

—Statistical modeling links cancer risk with number of stem cell divisions

Release Date: January 1, 2015
Addendum to news release added Jan. 7, 2015

Johns Hopkins Medicine is gratified by the responses and discussion generated by Cristian Tomasetti and Bert Vogelstein's research paper, "Variation in cancer risk among tissues can be explained by the number of stem cell divisions," published in Science on Jan. 2, 2015, and a news release describing the work, "Bad Luck of Random Mutations Plays Predominant Role in Cancer, Study Shows." Cancer is driven by a number of factors and causes, and concepts related to calculating risk are complex and often the subject of debate. To facilitate the ongoing discussion, and to address the many thoughtful questions their research stimulated, the two scientists have provided the following answers to frequently asked questions.

Is there an analogy that can help put the results of your research in perspective?

Getting cancer could be compared to getting into a car accident. Our results would be equivalent to showing a high correlation between length of trip and getting into an

FOR THE ME

Contacts:

Vanessa Wasta
410-614-2916
wasta@jhmi.edu

Download

Photo/Video files
this story are available
download.

We are providing this
understanding that it
help illustrate the story
corresponding news
find appropriate credit
images in the release
else you need, please
JHIMedia@jhmi.edu

View Files



Next Avenue @NextAvenue · Feb 9
Surprising Research on #Cancer and Bad Luck: is.gd/oFfIBj

Food&Drink Fanatic @foodanddrinkfan · Feb 8
Nope, cancer is not "just bad luck": Red meat contains saturated fats and the amino acid carnitine. If you con... cur.lv/treJ5

Cancer Daily @cancrddaily · Feb 8
Nope, cancer is not "just bad luck" - Arizona Daily Star bit.ly/1lz7oeR #cancer #health

Next Avenue @NextAvenue · Feb 8
If you're a #cancer survivor, you'll want to read this piece: is.gd/oFfIBj

Stem Cell News @mystemcellnews · Feb 8
Experts pan tying cancer to bad luck: JHU scientists Cristian Tomasetti and Bert Vogelstein had analysed stem... bit.ly/1KxLDCT



HEALTH CANCER

Most Types of Cancer Just 'Bad Luck,' Researchers Say

Helen Regan @hcregan | Jan. 2, 2015

Two thirds of cancers could be explained as biological misfortune

Researchers have found that bad luck plays a major role in determining most types of cancer, rather than genetics or risky lifestyle choices such as smoking.

The results, published in the journal Science on Thursday, found that random DNA mutations that amass in the body when stem cells divide into various tissues cause two thirds of cancers.

LYMPHOCYTES AND CANCER CELL

JUAN GARTNER—Getty Images/Science Photo Library RF

MOVING FORWARD

NEW QUESTIONS

Variation in cancer risk among tissues can be explained by the number of stem cell divisions

Cristian Tomasetti^{1*} and Bert Vogelstein^{2*}

Some tissue types give rise to human cancers millions of times more often than other tissue types. Although this has been recognized for more than a century, it has never been explained. Here, we show that the lifetime risk of cancers of many different types is strongly correlated (0.81) with the total number of divisions of the normal self-renewing cells maintaining that tissue's homeostasis. These results suggest that only a third of the variation in cancer risk among tissues is attributable to environmental factors or lifestyle.

Bad Luck of Random Mutations Plays Predominant Role in Cancer, Study Shows

—Statistical modeling links cancer risk with number of stem cell divisions

Release Date: January 1, 2015
Addendum to news release added Jan. 7, 2015

Johns Hopkins Medicine is gratified by the responses and discussion generated by Tomasetti's research paper, "Variation in cancer risk by the number of stem cell divisions," published in a news release describing the work, "Bad Luck of Random Mutations Plays Predominant Role in Cancer, Study Shows." Cancer is driven by random mutations, and concepts related to calculating risk are complex. To facilitate the ongoing discussion, and to address your research stimulated, the two scientists have prepared frequently asked questions.

How do you put the results of your research in perspective? For example, getting into a car accident. Our results would be that the relation between length of trip and getting into an

FOR THE MEDIA
Contacts:
Vanessa Wasta
410-614-2916
wasta@jhmi.edu

Download
Photo/Video files of this story are available for download.

We are providing this information to help you understand that it is not appropriate to use the images in the release for purposes other than those for which they were intended. If you need, please contact JHMedia@jhmi.edu

View Files

supplementary materials). With the number of clusters set equal to two, the tumors were classified in an unsupervised manner into one cluster with high ERS (9 tumor types) and another with low ERS (22 tumor types) (Fig. 2).

Next Avenue @NextAvenue - Feb 9
Surprising Research on #Cancer and Bad Luck: [s.g/d/9F9S](https://t.co/9F9S)

Food&Drink Fanatic @foodanddrinksfanatic - Feb 8
Nope, cancer is not "just bad luck": Fried meat contains saturated fats and the amino acid carnitine. If you con... [cur.Jw9e5](https://t.co/urJw9e5)

Cancer Daily @cancerdaily - Feb 8
Nope, cancer is not "just bad luck" [#health](#)

Next Avenue @NextAvenue - Feb 8
If you're a #cancer survivor, you

Stem Cell News @systemcellnews
Experts pin tying cancer to bad Bert Vogelstein had analysed st

HEALTH CANCER

Most Types of Cancer Just 'Bad Luck,' Researchers Say

Helen Regan @hcregan | Jan. 2, 2015

Two thirds of cancers could be explained as biological misfortune

After examining 31 cancer types, researchers found 22 were from mutations in stem cells that could not be prevented.

MOVING FORWARD

NEW QUESTIONS

Variation in cancer risk among tissues can be explained by the number of stem cell divisions

Cristian Tomasetti^{1*} and Bert Vogelstein^{2*}

Some tissue types give rise to human cancers millions of times more often than other tissue types. Although this has been recognized for more than a century, it has never been explained. Here, we show that the lifetime risk of cancers of many different types is strongly correlated (0.81) with the total number of divisions of the normal self-renewing cells maintaining that tissue's homeostasis. These results suggest that only a third of the variation in cancer risk among tissues is attributable to environmental factors or inherited predispositions. The majority is due to "bad luck," that is, random mutations arising during DNA replication in normal, noncancerous stem cells. This is important not only for understanding the disease but also for designing strategies to limit the mortality it causes.

Extrême variation in cancer incidence across different tissues is well known; for example, the lifetime risk of being diagnosed with cancer is 6.9% for lung, 1.08% for thyroid, 0.6% for brain and the rest of the nervous system, 0.003% for pelvic bone and 0.00072% for laryngeal cartilage (1-5). Some of these differences are associated with well-known risk factors such as smoking, alcohol use, ultraviolet light, or human papilloma virus (HPV) (4, 5), but this applies only to specific populations

exposed to potent mutagens or viruses. And such exposures cannot explain why cancer risk in tissues within the alimentary tract can differ by as much as a factor of 24 (esophagus (0.51%), large intestine (4.82%), small intestine (0.20%), and stomach (0.86%)) (3). Moreover, cancers of the small intestinal epithelium are three times less common than brain tumors (5), even though small intestinal epithelial cells are exposed to much higher levels of environmental mutagens than are cells within the brain, which are protected by the blood-brain barrier.

*Another well-studied contributor to cancer is



Bad Luck of Random Mutations Plays Predominant Role in Cancer, Study Shows

—Statistical modeling links cancer risk with number of stem cell divisions

Release Date: January 1, 2015
Addendum to news release added Jan. 7, 2015

Johns Hopkins Medicine is gratified by the responses and discussion generated by Cristian Tomasetti and Bert Vogelstein's research paper, "Variation in cancer risk among tissues can be explained by the number of stem cell divisions," published in Science on Jan. 2, 2015, and a news release describing the work, "Bad Luck of Random Mutations Plays Predominant Role in Cancer, Study Shows." Cancer is driven by a number of factors and causes, and concepts related to calculating risk are complex and often the subject of debate. To facilitate the ongoing discussion, and to address the many thoughtful questions their research stimulated, the two scientists have provided the following answers to frequently asked questions.

Is there an analogy that can help put the results of your research in perspective?

Getting cancer could be compared to getting into a car accident. Our results would be equivalent to showing a high correlation between length of trip and getting into an

FOR THE MEDIA

Contacts:

Vanessa Wasta
410-614-2916
wasta@jhmi.edu

Download

Photo/video files
this story are available
for download.

We are providing this understanding that it help illustrate the story corresponding news find appropriate credit images in the release else you need, please JHMedia@jhmi.edu

[View Files](#)

Wallace, **Paul**, Elhadad (2015) **What predicts media coverage of health science articles?** *AAAI Workshop on the World Wide Web and Public Health Intelligence.*

CS

MOVING FORWARD

SUMMARY

There's no end to exciting questions we can ask of big, open data

We need methods to link what people are saying on the web with real-world trends

This requires advancements at the intersection of language processing and data science



THANK YOU

WITH HELP FROM:

Advisors: Mark Dredze, Jason Eisner

Funding: Microsoft Research, NSF, JHU Dean's office

Flu:

-  David Broniatowski
-  Andrea Dugas
-  Nicholas Generous
-  Alex Lamb
-  Michael Smith




Medical search:

-  Eric Horvitz
-  Ryen White
-  Janice Tsai
-  Sara Javid
-  Luis Diaz

Air pollution:

-  Shiliang Wang
-  Angie Chen
-  Brian Schwartz

Doctor reviews:

-  Byron Wallace
-  Urmimala Sarkar
-  Thomas Trikalinos

Drug forums:

-  Meg Chisolm
-  Matthew Johnson
-  Ryan Vandrey

THANK YOU

QUESTIONS?

MOVING FORWARD

MORE DATA



The screenshot shows a vertical scroll of tweets. At the top, a tweet from Tiffany Thornton (@heresTiffany) says: "Dear Lord, please let whatever Chris has be food poisoning n not the flu. My birthday is Saturday and I really don't want to spend it sick." Below it, Roxeteraw (@RoxeteraRibbon) tweets: "On no!! Everyone at work had flu and now I think I've got :!!!!!! 🙏 please noooooooooooooo I have too many things to do!!". A tweet from Barack Obama follows, mentioning Shahid Kamal Ahmad (@shahidkamal) who says: "20 hour day, most of it work, one meal at 10:30pm, one toilet break in 13 hours and flu. And yet I'm ending the day totally psyched." Further down, Richard Oliver (@RichOliverActor) tweets: "Flu tabs taken & off to bed! leave you with another poem by Ricardo Pantelone as I head to my slumber #ActorLife x". The poem reads: "if i were a tree, what kind would i be? a mammoth oak, tall and slender? or a stumper version, silent, contemplative, tender? a weeping willow i would be, not for sadness, crying or misery, for the willows roots lie strong and deep, and by the winding rusby brook, in comfort sleep". Ricardo Pantelone is credited. Below the poem, a tweet from Jan Olsen and 2 others follows, mentioning NFD (@NFDvaccines) and saying: "If Teachers can #FightFlu by adding prevention messages to lessons. Ready-to-use work plans available bit.ly/1us6yK7 #K12". Then, Kamran M. Riaz (@kr156) tweets: "Can someone check if @bilimaher is down with the flu?As spokesperson for many atheists, his silence belies his approval #ChapelHillShooting". Next, Kathleen Bachynski and 10 others follow, mentioning CIDRAP (@CIDRAP) and saying: "FLU SCAN: Parotitis in flu patients; Global flu update; H7N9 in China; Avian flu in Taiwan, Bulgaria ow.ly/1Uj0z2 #G". Below that, NFD (@NFDvaccines) tweets: "Don't weather the #flu! When flu hits, act fast! #FightFlu ow.ly/1M3z4". A photo of a stormy sky with lightning is shown. Then, Syndromic.org (@ISDS) tweets: "HK's Dr. Ko Wing-man on Flu Reassessment Concerns (Avian Flu Diary) - bit.ly/1MASWc". Below that, alex vespignani and 11 others follow, mentioning Skoll Global (@SkollGlobal) and saying: "Flu Near You featured on Fighting the Flu - FOX 8 WWUE New Orleans fox8live.com/Clip/1125348/... #FluNearYou". Then, NFD (@NFDvaccines) tweets: "Prompt use of antivirals is key this #flu season via @CDCFlu ow.ly/1Z24G". Finally, Syndromic.org (@ISDS) tweets: "US Flu Activity Down Slightly, but Elderly Hit Hard (CIDRAP) - bit.ly/1M83w3q". At the bottom, a tweet from Elin Rasmussen (@ElinRasmussen) says: "First Flu season under winter flu, but claims to be div.it:8VYNbc".

“What opinions are people tweeting about gun rights?”

Sprite can model this!

We created a structured prior that incorporates geographic data about gun ownership

- Certain topics have a higher/lower prior depending on whether a tweet is from a high/low gun ownership state

Alabama	2,623	1,294	51.7	1,329	48.3
Alaska	2,716	1,627	57.8	1,089	42.2
Arizona	3,066	989	31.1	2,077	68.9
Arkansas	2,780	1,431	55.3	1,349	44.7
California*	3,897	846	21.3	3,051	78.7
Colorado	1,947	629	34.7	1,318	65.3
Connecticut*	7,449	1,279	16.7	6,170	83.3
Delaware*	3,421	934	25.5	2,487	74.5
The District	1,859	66	3.8	1,793	96.2
Florida*	4,454	1,072	24.5	3,382	75.5
Georgia	4,277	1,745	40.3	2,532	59.7
Hawaii*	4,450	477	8.7	3,973	91.3
Idaho	4,430	2,394	55.3	2,036	44.7
Illinois*	2,103	396	20.2	1,707	79.8

Benton, Paul, Hancock, Dredze.
A structured model of topic and perspective in social media.
Submitted to *ICWSM*.

MOVING FORWARD

MORE DATA

Sprite can model this!

violence
culture
less
sense
problem
makes
world
country
common
america's

children
kids
likely
times
murdered
giving
nothing
stand
sorry
insane

nra
war
even
keep
murder
members
liberal
fear
government
call

teachers
school
armed
schools
carry
god
kids
teacher
protect
security

Associated with: Low gun ownership

High gun ownership

MOVING FORWARD

MORE DATA

This topic is associated with high gun ownership:

guns
'merica
truck
shoot
deer
hunting
day
beer
season
friends

This topic is associated with high gun ownership:

guns
'merica
truck
shoot
deer
hunting
day
beer
season
friends

Probably also associated with:

- Population density
- Political affiliation
- ...

MOVING FORWARD



Prior for triple (i,j,k) :

$$\tilde{\phi}_{(i,j,k)v} = \exp(\omega_{iv}^{(\text{drug})} + \omega_{jv}^{(\text{route})} + \omega_{kv}^{(\text{aspect})})$$

$$\phi_{(i,j,k)} \sim \text{Dirichlet}(\tilde{\phi}_{(i,j,k)})$$



distribution over words for this triple

In general, prior for tuple \mathbf{t} :

$$\tilde{\phi}_{\vec{t}v} = \exp\left(\sum_{k=1}^K \omega_{\vec{t}_k v}^{(k)}\right)$$

$$\phi_{\vec{t}} \sim \text{Dirichlet}(\tilde{\phi}_{\vec{t}})$$

SPRITE

DEFINITION

β_i has one value of 1 in each factor; 0 elsewhere

$$\tilde{\phi}_{iv} = \exp\left(\sum_{c=1}^{C(\phi)} \beta_{ic} \omega_{cv}\right)$$

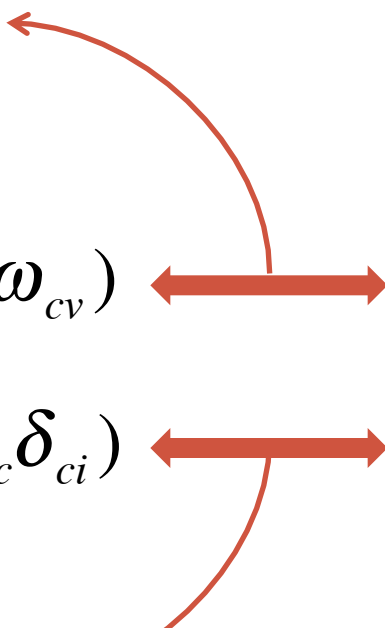
$$\tilde{\theta}_{mi} = \exp\left(\sum_{c=1}^{C(\theta)} \alpha_{mc} \delta_{ci}\right)$$

δ_c is transpose of β_i

Recall (FLDA):

$$\tilde{\phi}_{\vec{t}v} = \exp\left(\sum_{k=1}^K \omega_{\vec{t}_k v}^{(k)}\right)$$

$$\tilde{\theta}_{m\vec{t}} = \exp\left(\sum_{k=1}^K \alpha_{m\vec{t}_k}^{(k)}\right)$$



SPRITE

DEFINITION

β_i has one value of 1 in each factor; 0 elsewhere

$$\tilde{\phi}_{iv} = \exp\left(\sum_{c=1}^{C(\phi)} \beta_{ic} \omega_{cv}\right)$$

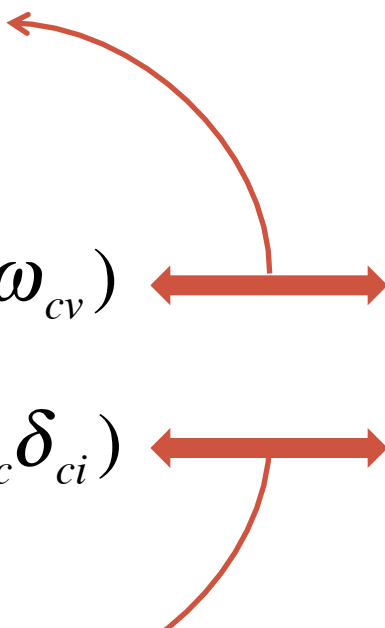
$$\tilde{\theta}_{mi} = \exp\left(\sum_{c=1}^{C(\theta)} \alpha_{mc} \delta_{ci}\right)$$

δ_c is transpose of β_i

Recall (FLDA):

$$\tilde{\phi}_{\vec{t}v} = \exp\left(\sum_{k=1}^K \omega_{\vec{t}_k v}^{(k)}\right)$$

$$\tilde{\theta}_{m\vec{t}} = \exp\left(\sum_{k=1}^K \alpha_{m\vec{t}_k}^{(k)}\right)$$



Suppose we want to model how **perspective** influences topics

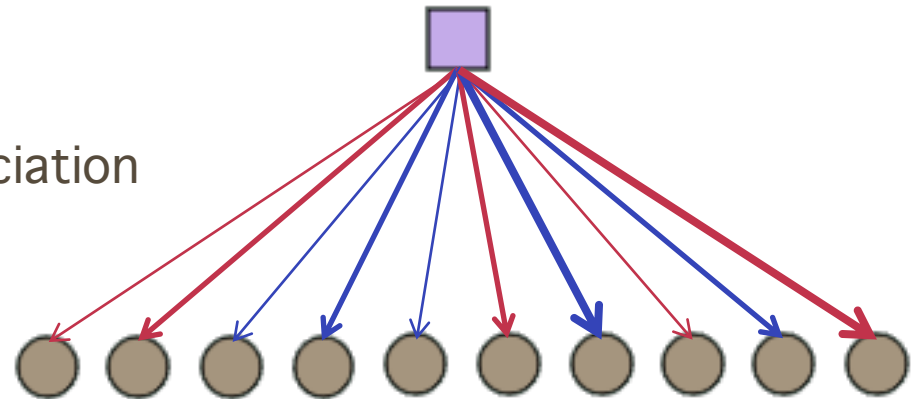
- e.g. certain topics are “**pro**” or “**anti**” gun control

A single-component SPRITE model:

$$\tilde{\phi}_{kv} = \exp(r_k \omega_v)$$

The v th word's perspective association

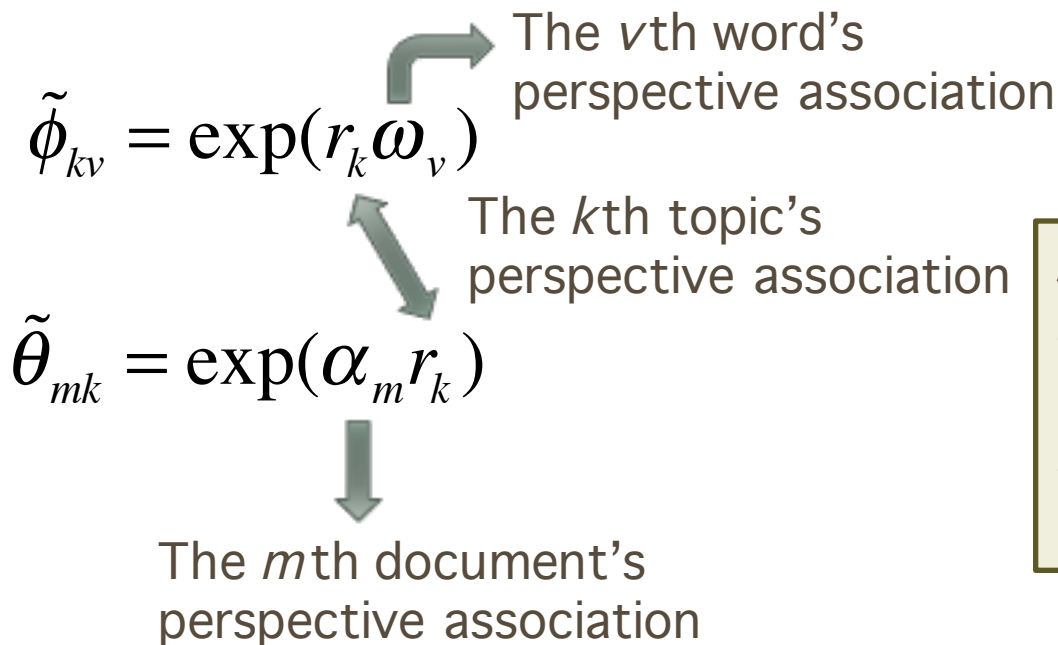
The k th topic's perspective association



Suppose we want to model how **perspective** influences topics

- e.g. certain topics are “**pro**” or “**anti**” gun control

A single-component SPRITE model:



A positive r_k means:


- Words with positive ω_k are more likely in topic k
- Topic k is more likely in documents with positive α_m

Suppose we want to model how **perspective** influences topics

- e.g. certain topics are “**pro**” or “**anti**” gun control

A single-component SPRITE model:

Incorporating soft supervision: $\alpha_m \sim \mathcal{N}(s_m, \sigma^2)$

$$\tilde{\theta}_{mk} = \exp(\alpha_m r_k)$$


Supervision s_m is a function of:

- Survey data (% gun ownership in each state)
- Hashtags (#GunControlNow vs #NoGunControl)

MOVING FORWARD



MOVING FORWARD

BETTER DATA

Previous section: linking text to **population data**

Another idea: linking text to **individual data**

MOVING FORWARD

BETTER DATA

Previous section: linking text to population data

Another idea: linking text to individual data

Extraversion:



Introversion:



Schwartz et al. (2013) **Personality, gender, and age in the language of social media: the open-vocabulary approach.** *PLOS ONE*.

MOVING FORWARD

BETTER DATA

Previous section: linking text to population data

Another idea: linking text to individual data

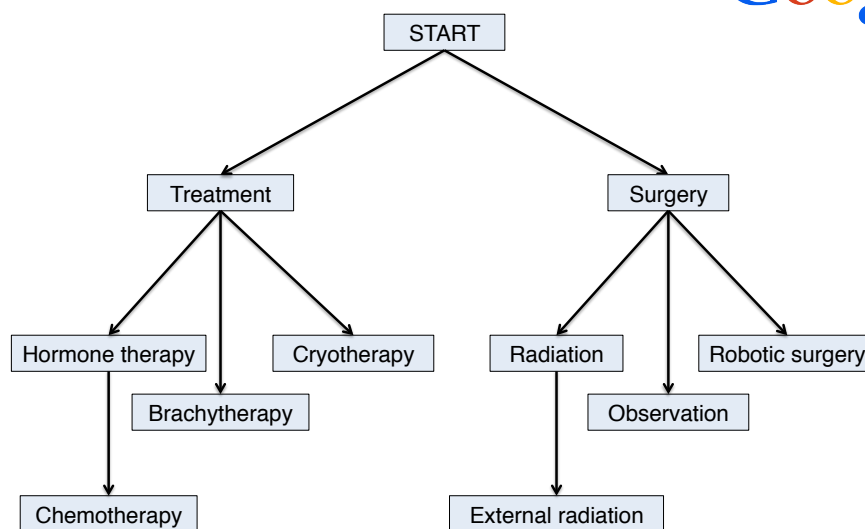
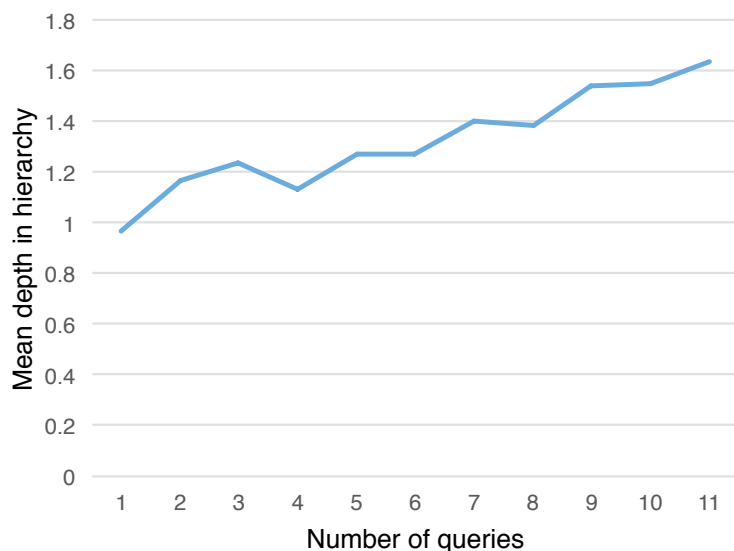
- New territory: **clinical records**



MOVING FORWARD

BETTER DATA

How does the web influence medical decision-making?



Paul, White, Horvitz (2015) **Web search as medical decision support for cancer.** *International World Wide Web Conference (WWW).*

