~~Anchored Correlation~~

# How to Topic Model with Literally Thousands of Information Bottlenecks 🍾

~~an Knowledge~~

## Ryan J. Gallagher

🐦 @ryanjgallag

github.com/gregversteeg/corex_topic

USC University of Southern California
*Information Sciences Institute*

netsi

Northeastern University
*Network Science Institute*

VACC
VERMONT ADVANCED COMPUTING CORE
UNIVERSITY OF VERMONT

VERMONT COMPLEX SYSTEMS CENTER
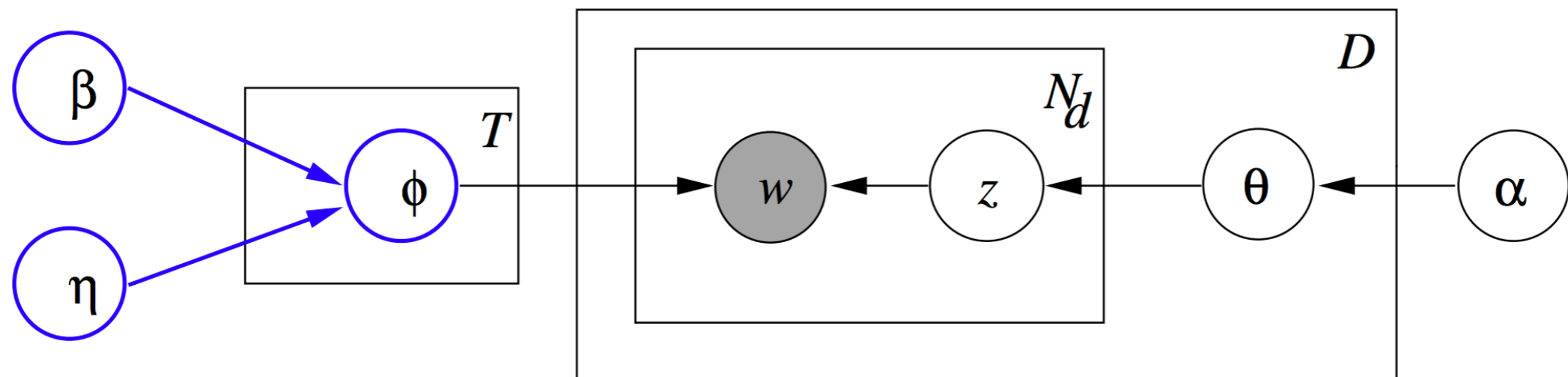
# LDA is a *generative* topic model

# LDA is a *generative* topic model
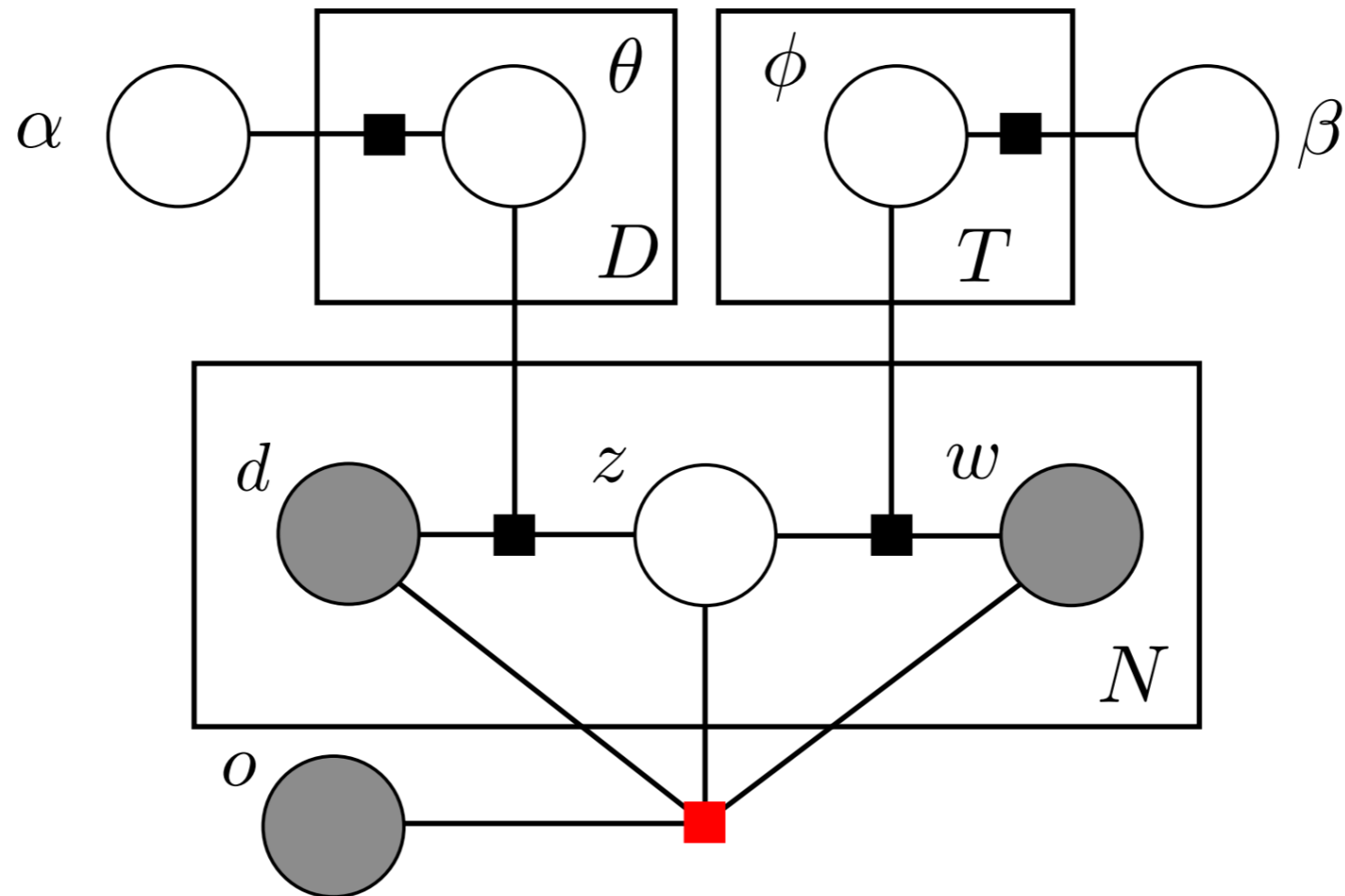
## The Good:

Priors explicitly encode your beliefs about what topics can be, and easily allow for iterative development of new topic models
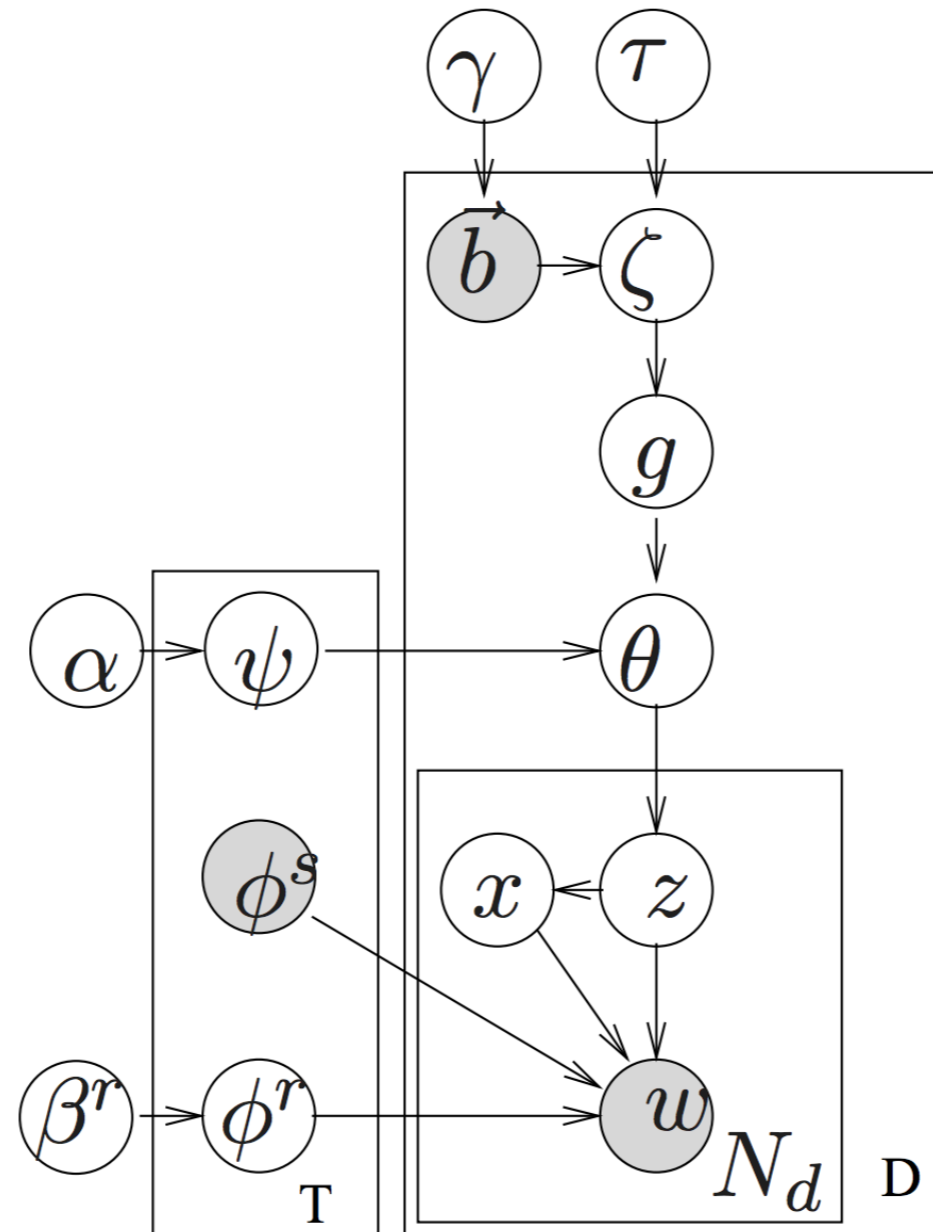
# Domain Knowledge via Dirichlet Forest Priors



"Incorporating Domain Knowledge into Topic Modeling via Dirichlet Forest Priors." Andrzejewski et al. *ICML* (2009)

# Domain Knowledge via First-Order Logic



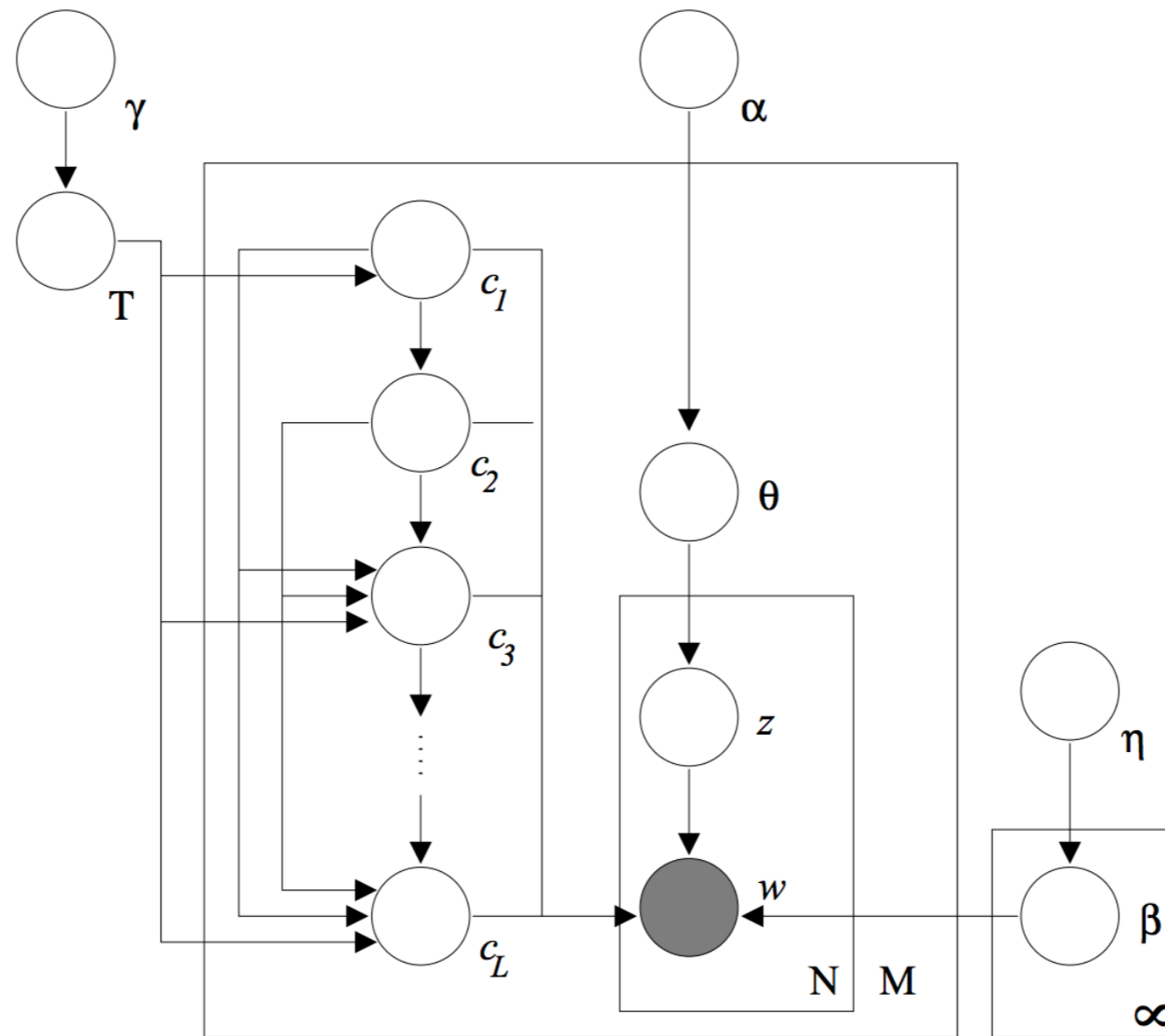"A Framework for Incorporating General Domain Knowledge into Latent Dirichlet Allocation Using First-Order Logic." Andrzejewski et al. *IJCAI* (2011).

# SeededLDA



"Incorporating Lexical Priors into Topic Models." Jagarlamudi et al. *EACL* (2012)

# Hierarchical LDA



"Hierarchical Topic Models and the Nested Chinese Restaurant Process." Griffiths et al. *Neural Information Processing Systems* (2003).

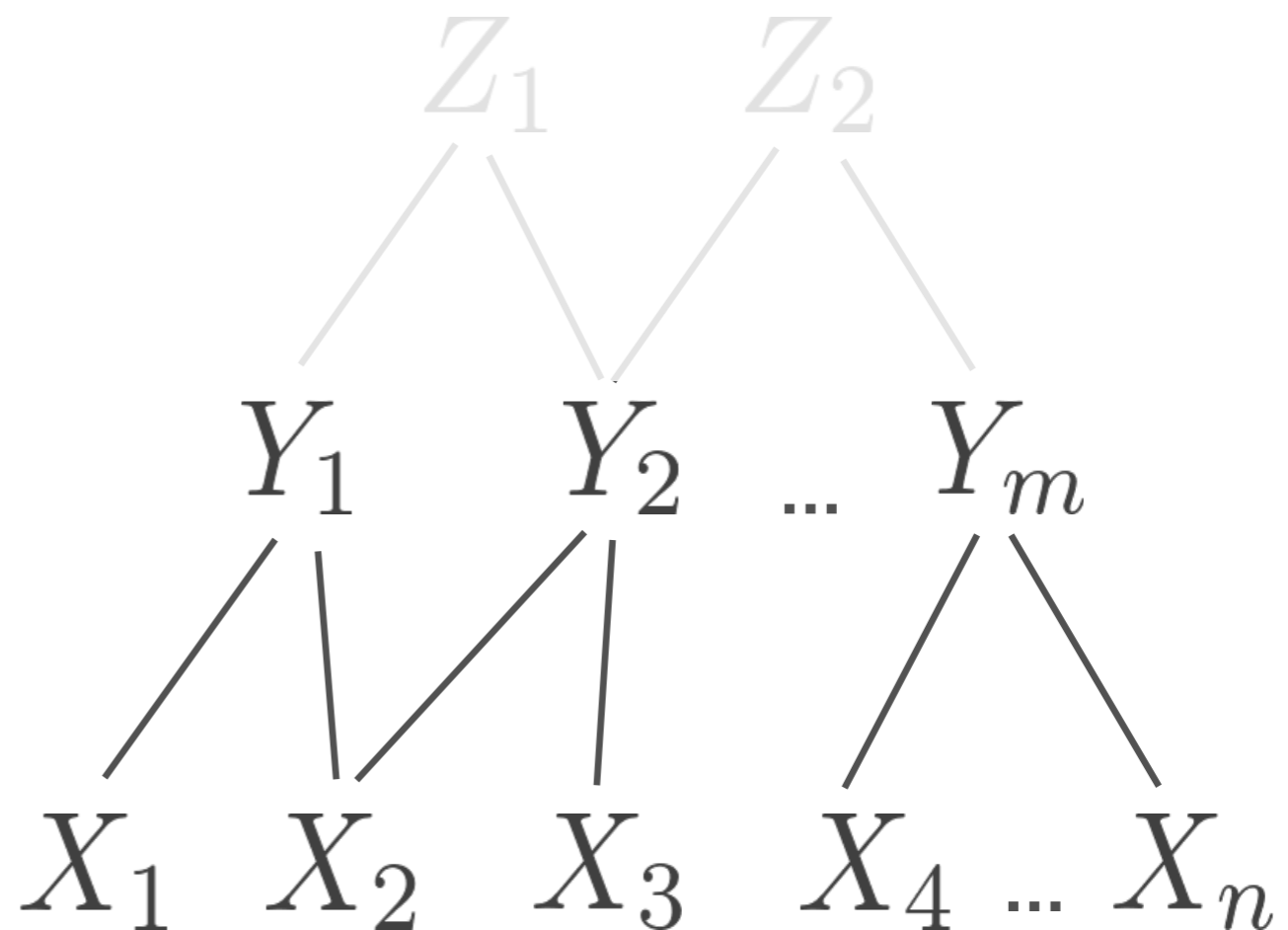# A Generative Modeling Tradeoff

## The Good:

Priors explicitly encode your beliefs about what topics can be, and easily allow for iterative development of new topic models

## The Bad:

Each additional prior takes a very specific view of the problem at hand, which both limits what a topic can be and makes it harder to justify in applications and to domain experts

We propose a topic model that learns topics through information-theoretic criteria, rather than a generative model

$$Z_1 \qquad Z_2$$

$$Y_1 \qquad Y_2 \quad ... \quad Y_m$$

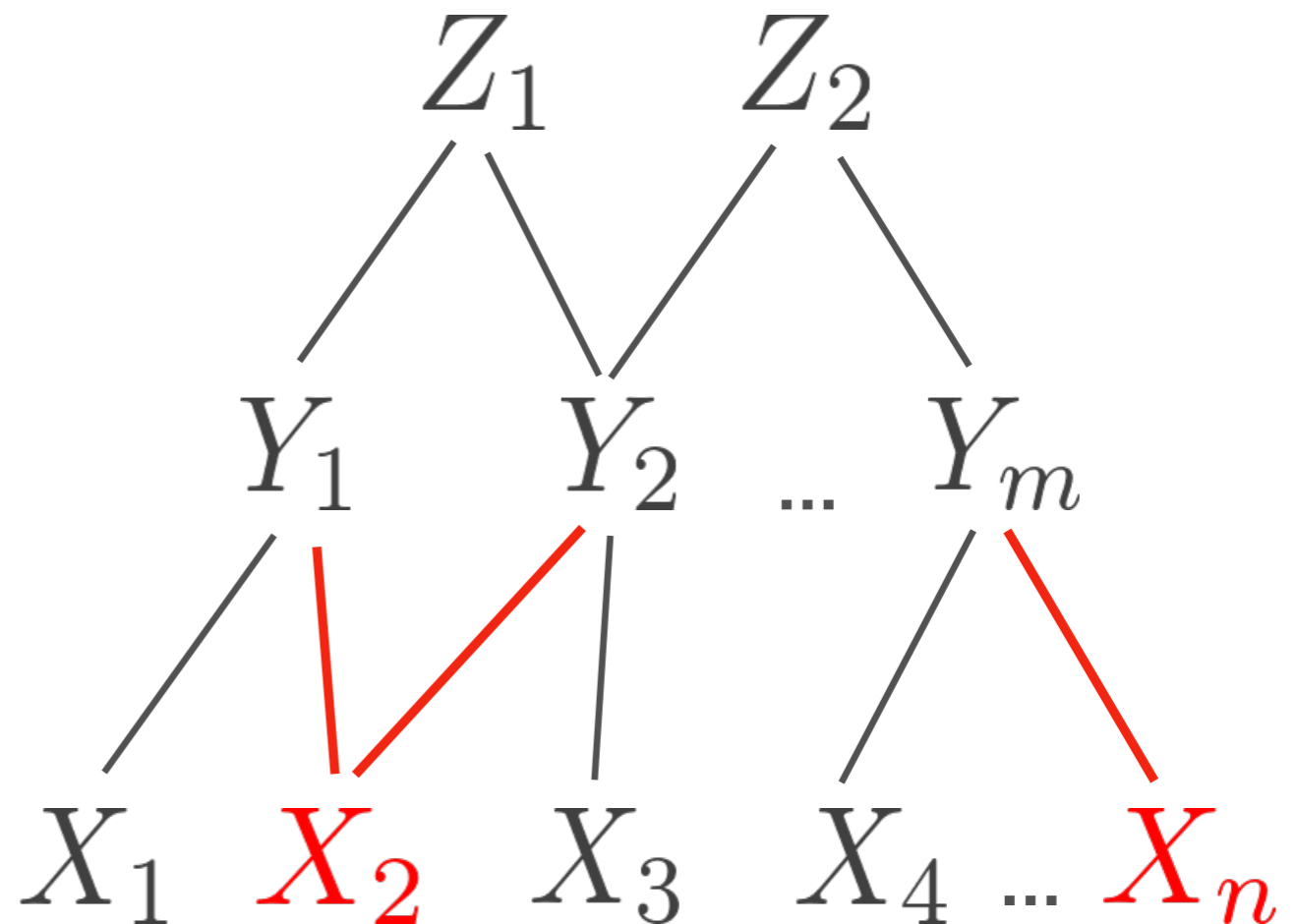$$X_1 \quad X_2 \quad X_3 \quad X_4 \quad ... \quad X_n$$

@ryanjgallag

# Proposed Work

We propose a topic model that learns topics through information-theoretic criteria, rather than a generative model, within a framework that yields **hierarchical** and **semi-supervised** extensions with *no additional assumptions*

@ryanjgallag

# A Different Perspective on "Topics"

Consider three documents:

$$d_1 \qquad\qquad d_2 \qquad\qquad d_3$$

| $x_1 \quad x_2$ | $x_3 \quad x_4$ | $x_5$ |

$$(1,1,0,0,0) \qquad (0,0,1,1,0) \qquad (0,0,0,0,1)$$

**LDA:** a topic is a distribution over words



$$P(Y=1)=1 \xleftarrow{\quad d_1 \qquad d_3 \qquad d_2 \quad} P(Y=2)=1$$

# A Different Perspective on "Topics"

Consider three documents:

$$d_1 \qquad\qquad d_2 \qquad\qquad d_3$$

| $x_1 \quad x_2$ |
| :---: |

| $x_3 \quad x_4$ |
| :---: |

| $x_5$ |
| :---: |

$$(1, 1, 0, 0, 0) \qquad (0, 0, 1, 1, 0) \qquad (0, 0, 0, 0, 1)$$

**LDA:** a topic is a distribution over words          **CorEx:** a topic is a binary latent factor

Consider three documents:

$$d_1$$

| $x_1$ | $x_2$ |
|-------|-------|

$(1, 1, 0, 0, 0)$

$$d_2$$

| $x_3$ | $x_4$ |
|-------|-------|

$(0, 0, 1, 1, 0)$

$$d_3$$

| $x_5$ |
|-------|

$(0, 0, 0, 0, 1)$

**LDA:** a topic is a distribution over words

**CorEx:** a topic is a binary latent factor



$$P(Y = 1) = 1 \xleftarrow{d_1 \quad d_3 \quad d_2} P(Y = 2) = 1$$

$Y_1 \qquad Y_2$

$x_1 \quad x_2 \qquad x_3 \quad x_4 \quad x_5$

$$P(Y_1 = 1) \begin{array}{l} d_1 \\ \\ d_3 \quad d_2 \end{array}$$

$$P(Y_2 = 1)$$

# CorEx Objective (example)

Documents

Probability table

$$d_1$$

| $x_1$ | $x_2$ |
|---|---|

$$(1, 1, 0, 0, 0)$$

$$d_2$$

| $x_3$ | $x_4$ |
|---|---|

$$(0, 0, 1, 1, 0)$$

|  | $X_2 = 0$ | $X_2 = 1$ |
|---|---|---|
| $X_1 = 0$ | 1/2 | 0 |
| $X_1 = 1$ | 0 | 1/2 |

Documents

Probability table

$$d_1 \qquad\qquad d_2$$

| $x_1$ | $x_2$ |
|-------|-------|

| $x_3$ | $x_4$ |
|-------|-------|

$$(1,1,0,0,0) \qquad (0,0,1,1,0)$$

|  | $X_2 = 0$ | $X_2 = 1$ |
|---------|-----------|-----------|
| $X_1 = 0$ | 1/2 | 0 |
| $X_1 = 1$ | 0 | 1/2 |

Words 1 and 2 are related:

$$I(X_1 : X_2) = D_{KL}\big(p(x_1, x_2) \,||\, p(x_1)p(x_2)\big) = 1 \text{ bit}$$

# CorEx Objective (example)

Documents

Probability table

$$d_1 \qquad d_2$$

| $x_1$ | $x_2$ |
|---|---|

| $x_3$ | $x_4$ |
|---|---|

$(1,1,0,0,0) \qquad (0,0,1,1,0)$

|  | $X_2 = 0$ | $X_2 = 1$ |
|---|---|---|
| $X_1 = 0$ | 1/2 | 0 |
| $X_1 = 1$ | 0 | 1/2 |

Words 1 and 2 are related:

$$I(X_1 : X_2) = D_{KL}\big(p(x_1, x_2) \,||\, p(x_1)p(x_2)\big) = 1 \text{ bit}$$

Hypothesize a latent factor: $Y_1 = X_1 = X_2$

$$Y_1$$

$$X_1 \qquad X_2$$

Documents

Probability table

$$d_1$$

| $x_1$ | $x_2$ |
|---|---|

$$(1, 1, 0, 0, 0)$$

$$d_2$$

| $x_3$ | $x_4$ |
|---|---|

$$(0, 0, 1, 1, 0)$$

|  | $X_2 = 0$ | $X_2 = 1$ |
|---|---|---|
| $X_1 = 0$ | 1/2 | 0 |
| $X_1 = 1$ | 0 | 1/2 |

Words 1 and 2 are related:

$$I(X_1 : X_2) = D_{KL}\big(p(x_1, x_2) \,||\, p(x_1)p(x_2)\big) = 1 \text{ bit}$$

$$Y_1$$

$$X_1 \qquad X_2$$

Hypothesize a latent factor: $Y_1 = X_1 = X_2$

Then conditioned on $Y_1$, words 1 and 2 are independent

$$D_{KL}\big(p(x_1, x_2 \mid y_1) \,||\, p(x_1 \mid y_1)p(x_2 \mid y_1)\big) = 0 \text{ bits}$$

# CorEx Objective (example)

Documents

$$d_1 \qquad d_2$$

| $x_1$ $x_2$ | | $x_3$ $x_4$ |
| --- | --- | --- |

$(1, 1, 0, 0, 0) \qquad (0, 0, 1, 1, 0)$

Probability table

| | $X_2 = 0$ | $X_2 = 1$ |
| --- | --- | --- |
| $X_1 = 0$ | 1/2 | 0 |
| $X_1 = 1$ | 0 | 1/2 |

Words 1 and 2 are related:

$$I(X_1 : X_2) = D_{KL}\big(p(x_1, x_2) \,||\, p(x_1)p(x_2)\big) = 1 \text{ bit}$$

$Y_1$

Hypothesize a latent factor: $Y_1 = X_1 = X_2$

$X_1 \qquad X_2$

Then conditioned on $Y_1$, words 1 and 2 are independent

$$D_{KL}\big(p(x_1, x_2 \mid y_1) \,||\, p(x_1 \mid y_1)p(x_2 \mid y_1)\big) = 0 \text{ bits}$$

**Goal:** find latent factors that make words conditionally independent

# CorEx Objective

**Goal:** find latent factors that make words conditionally independent

$$\min_{Y} D_{KL}\left( p(x_1, x_2, \dots x_n \mid y) \,\|\, \prod_i p(x_i \mid y) \right)$$

# CorEx Objective

**Goal:** find latent factors that make words conditionally independent

$$\min_Y D_{KL}\left( p(x_1, x_2, \dots x_n \mid y) \,||\, \prod_i p(x_i \mid y) \right) = \min_Y \underline{TC(X_1, X_2, \dots, X_N \mid Y)}$$

Total correlation conditioned on *Y*

# CorEx Objective

**Goal:** find latent factors that make words conditionally independent



$$\min_Y D_{KL}\left(p(x_1, x_2, \ldots x_n \mid y) \,||\, \prod_i p(x_i \mid y)\right) = \min_Y TC(X_1, X_2, \ldots, X_N \mid Y)$$

$TC(X \mid Y) = 0$ if and only if the topic "explains" all the dependencies (total correlation)

Hence, "Total **Cor**relation **Ex**planation" (CorEx)

# CorEx Objective

**Goal:** find latent factors that make words conditionally independent

$$Y_1 \qquad Y_2$$

$$X_1 \qquad X_2 \qquad X_3 \qquad X_4 \qquad X_5$$

$$\min_Y D_{KL}\left( p(x_1, x_2, \ldots x_n \mid y) \,\|\, \prod_i p(x_i \mid y) \right) = \min_Y TC(X_1, X_2, \ldots, X_N \mid Y)$$

In order to maximize the information $TC(X_{G_j})$ between a group of words $G_j$ in topic $j$ we consider a tractable lower bound:

$$TC(X_{G_j}) - TC(X_{G_j} \mid Y_j) \leq TC(X_{G_j})$$
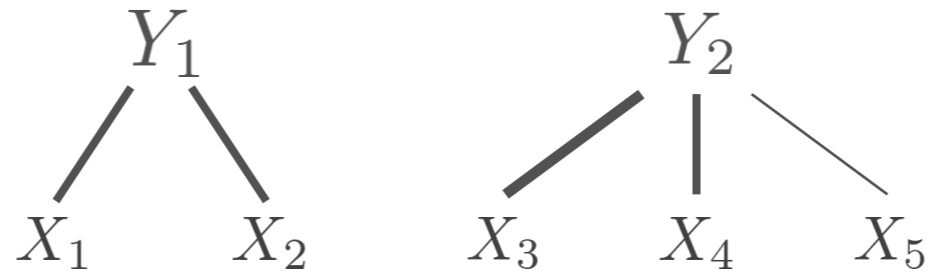
@ryanjgallag

# CorEx Objective

**Goal:** find latent factors that make words conditionally independent

$$\min_Y D_{KL}\left( p(x_1, x_2, \ldots x_n \mid y) \,\|\, \prod_i p(x_i \mid y) \right) = \min_Y TC(X_1, X_2, \ldots, X_N \mid Y)$$

In order to maximize the information $TC(X_{G_j})$ between a group of words $G_j$ in topic $j$ we consider a tractable lower bound:

$$TC(X_{G_j}) - TC(X_{G_j} \mid Y_j) \leq TC(X_{G_j})$$

We maximize this lower bound over $m$ topics

$$\max_{G_j, p(y_j \mid x_{G_j})} \sum_{j=1}^{m} TC(X_{G_j}) - TC(X_{G_j} \mid Y_j)$$

# CorEx Objective

We can now rewrite the objective:

$$\max_{G_j, p(y_j | x_{G_j})} \sum_{j=1}^{m} TC(X_{G_j}) - TC(X_{G_j} \mid Y_j) = \max_{G_j, p(y_j | x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

# CorEx Objective

We can now rewrite the objective:

$$\max_{G_j, p(y_j|x_{G_j})} \sum_{j=1}^{m} TC(X_{G_j}) - TC(X_{G_j} \mid Y_j) = \max_{G_j, p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

We transform this from a combinatorial to a continuous optimization by introducing variables $\alpha_{i,j} \in [0, 1]$ and "relaxing" words into informative topics

$$\max_{G_j, p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} \alpha_{i,j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

# CorEx Objective

We can now rewrite the objective:

$$\max_{G_j,p(y_j|x_{G_j})} \sum_{j=1}^{m} TC(X_{G_j}) - TC(X_{G_j} \mid Y_j) = \max_{G_j,p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

We transform this from a combinatorial to a continuous optimization by introducing variables $\alpha_{i,j} \in [0,1]$ and "relaxing" words into informative topics

$$\max_{G_j,p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} \alpha_{i,j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

This relaxation yields a set of update equations which we can iterate through until convergence

We can now rewrite the objective:

$$\max_{G_j,p(y_j|x_{G_j})} \sum_{j=1}^{m} TC(X_{G_j}) - TC(X_{G_j} \mid Y_j) = \max_{G_j,p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i\in G_j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

We transform this from a combinatorial to a continuous optimization by introducing variables $\alpha_{i,j} \in [0,1]$ and "relaxing" words into informative topics

$$\max_{G_j,p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i\in G_j} \alpha_{i,j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

**Under the hood:**
1. We introduce a sparsity optimization for the update equations,
$$O(N_{\text{docs}} n_{\text{types}}) \to O(N_{\text{docs}}) + O(n_{\text{types}}) + O(\rho_{\text{tokens}})$$

# CorEx Objective

We can now rewrite the objective:

$$\max_{G_j, p(y_j | x_{G_j})} \sum_{j=1}^{m} TC(X_{G_j}) - TC(X_{G_j} \mid Y_j) = \max_{G_j, p(y_j | x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

We transform this from a combinatorial to a continuous optimization by introducing variables $\alpha_{i,j} \in [0,1]$ and "relaxing" words into informative topics

$$\max_{G_j, p(y_j | x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} \alpha_{i,j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

**Under the hood:**
1.  We introduce a sparsity optimization for the update equations,

$$O(N_{\mathrm{docs}} n_{\mathrm{types}}) \to O(N_{\mathrm{docs}}) + O(n_{\mathrm{types}}) + O(\rho_{\mathrm{tokens}})$$

   by assuming words are represented by binary random variables

# CorEx Objective

We can now rewrite the objective:

$$\max_{G_j,\,p(y_j|x_{G_j})} \sum_{j=1}^{m} TC(X_{G_j}) - TC(X_{G_j} \mid Y_j) = \max_{G_j,\,p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

We transform this from a combinatorial to a continuous optimization by introducing variables $\alpha_{i,j} \in [0,1]$ and "relaxing" words into informative topics

$$\max_{G_j,\,p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} \alpha_{i,j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

**Under the hood:**
1. We introduce a sparsity optimization for the update equations,

   $$O(N_{\mathrm{docs}} n_{\mathrm{types}}) \to O(N_{\mathrm{docs}}) + O(n_{\mathrm{types}}) + O(\rho_{\mathrm{tokens}})$$

   by assuming words are represented by binary random variables

2. The current relaxation scheme places each word in one topic, resulting in a partition of the vocabulary, rather than mixed membership topics

# CorEx Objective

We can now rewrite the objective:

$$\max_{G_j,p(y_j|x_{G_j})} \sum_{j=1}^{m} TC(X_{G_j}) - TC(X_{G_j} \mid Y_j) = \max_{G_j,p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

We transform this from a combinatorial to a continuous optimization by introducing variables $\alpha_{i,j} \in [0,1]$ and "relaxing" words into informative topics

$$\max_{G_j,p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} \alpha_{i,j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

**Under the hood:**
1. We introduce a sparsity optimization for the update equations,
   $$O(N_{\text{docs}} n_{\text{types}}) \to O(N_{\text{docs}}) + O(n_{\text{types}}) + O(\rho_{\text{tokens}})$$
   by assuming words are represented by binary random variables

2. The current relaxation scheme places each word in one topic, resulting in a partition of the vocabulary, rather than mixed membership topics

These are issues of speed, not theory

# CorEx Topic Examples

**Data:** news articles about Hillary Clinton's presidential campaign, up to August 2016

Work by Abigail Ross and the Computational Story Lab, University of Vermont

@ryanjgallag

# CorEx Topic Examples

**Data:** news articles about Hillary Clinton's presidential campaign, up to August 2016

## Clinton Article Topics

**1:** server, department, classified, information, private, investigation, fbi, email, emails, secretary

**3:** sanders, bernie, primary, vermont, win, voters, race, nomination, vote, polls

**6:** crowd, woman, speech, night, women, stage, man, mother, audience, life

**8:** percent, poll, points, percentage, margin, survey, according, 10, polling, university

**9:** federal, its, officials, law, including, committee, staff, statement, director, group

**13:** islamic, foreign, military, terrorism, war, syria, iraq, isis, u, terrorist

**14:** trump, donald, trump's, republican, nominee, party, convention, top, election, him

Work by Abigail Ross and the Computational Story Lab, University of Vermont

# CorEx Topic Examples

**Data:** news articles about Hillary Clinton's presidential campaign, up to August 2016

## Clinton Article Topics

**1:** server, department, classified, information, private, investigation, fbi, email, emails, secretary

Words ranked by mutual information with topic

**3:** sanders, bernie, primary, vermont, win, voters, race, nomination, vote, polls

**6:** crowd, woman, speech, night, women, stage, man, mother, audience, life

**8:** percent, poll, points, percentage, margin, survey, according, 10, polling, university

**9:** federal, its, officials, law, including, committee, staff, statement, director, group

**13:** islamic, foreign, military, terrorism, war, syria, iraq, isis, u, terrorist

**14:** trump, donald, trump's, republican, nominee, party, convention, top, election, him

Work by Abigail Ross and the Computational Story Lab, University of Vermont

# CorEx Topic Examples

**Data:** news articles about Hillary Clinton's presidential campaign, up to August 2016

Topics ranked
by total
correlation

## Clinton Article Topics

**1:** server, department, classified, information, private, investigation, fbi, email, emails, secretary

**3:** sanders, bernie, primary, vermont, win, voters, race, nomination, vote, polls

**6:** crowd, woman, speech, night, women, stage, man, mother, audience, life

**8:** percent, poll, points, percentage, margin, survey, according, 10, polling, university

**9:** federal, its, officials, law, including, committee, staff, statement, director, group

**13:** islamic, foreign, military, terrorism, war, syria, iraq, isis, u, terrorist

**14:** trump, donald, trump's, republican, nominee, party, convention, top, election, him

Work by Abigail Ross and the Computational Story Lab, University of Vermont

# CorEx Topic Examples

**Data:** news articles about Hillary Clinton's presidential campaign, up to August 2016

## Clinton Article Topics

Most informative topic

**1:** server, department, classified, information, private, investigation, fbi, email, emails, secretary

**3:** sanders, bernie, primary, vermont, win, voters, race, nomination, vote, polls

**6:** crowd, woman, speech, night, women, stage, man, mother, audience, life

**8:** percent, poll, points, percentage, margin, survey, according, 10, polling, university

**9:** federal, its, officials, law, including, committee, staff, statement, director, group

**13:** islamic, foreign, military, terrorism, war, syria, iraq, isis, u, terrorist

**14:** trump, donald, trump's, republican, nominee, party, convention, top, election, him

Work by Abigail Ross and the Computational Story Lab, University of Vermont

# CorEx Topic Examples

**Data:** news articles about Hillary Clinton's presidential campaign, up to August 2016

## Clinton Article Topics

**1:** server, department, classified, information, private, investigation, fbi, email, emails, secretary

**3:** sanders, bernie, primary, vermont, win, voters, race, nomination, vote, polls

**6:** crowd, woman, speech, night, women, stage, man, mother, audience, life

**8:** percent, poll, points, percentage, margin, survey, according, 10, polling, university

**9:** federal, its, officials, law, including, committee, staff, statement, director, group

**13:** islamic, foreign, military, terrorism, war, syria, iraq, isis, u, terrorist

**14:** trump, donald, trump's, republican, nominee, party, convention, top, election, him

Work by Abigail Ross and the Computational Story Lab, University of Vermont

# CorEx Topic Examples

**Data:** news articles about Hillary Clinton's presidential campaign, up to August 2016

## Clinton Article Topics

**1:** server, department, classified, information, private, investigation, fbi, email, emails, secretary

**3:** sanders, bernie, primary, vermont, win, voters, race, nomination, vote, polls

**6:** crowd, woman, speech, night, women, stage, man, mother, audience, life

**8:** percent, poll, points, percentage, margin, survey, according, 10, polling, university

**9:** federal, its, officials, law, including, committee, staff, statement, director, group

**13:** islamic, foreign, military, terrorism, war, syria, iraq, isis, u, terrorist

**14:** trump, donald, trump's, republican, nominee, party, convention, top, election, him

Work by Abigail Ross and the Computational Story Lab, University of Vermont

# CorEx Topic Examples

**Data:** news articles about Hillary Clinton's presidential campaign, up to August 2016

## Clinton Article Topics

**1:** server, department, classified, information, private, investigation, fbi, email, emails, secretary

**3:** sanders, bernie, primary, vermont, win, voters, race, nomination, vote, polls

**6:** crowd, woman, speech, night, women, stage, man, mother, audience, life

**8:** percent, poll, points, percentage, margin, survey, according, 10, polling, university
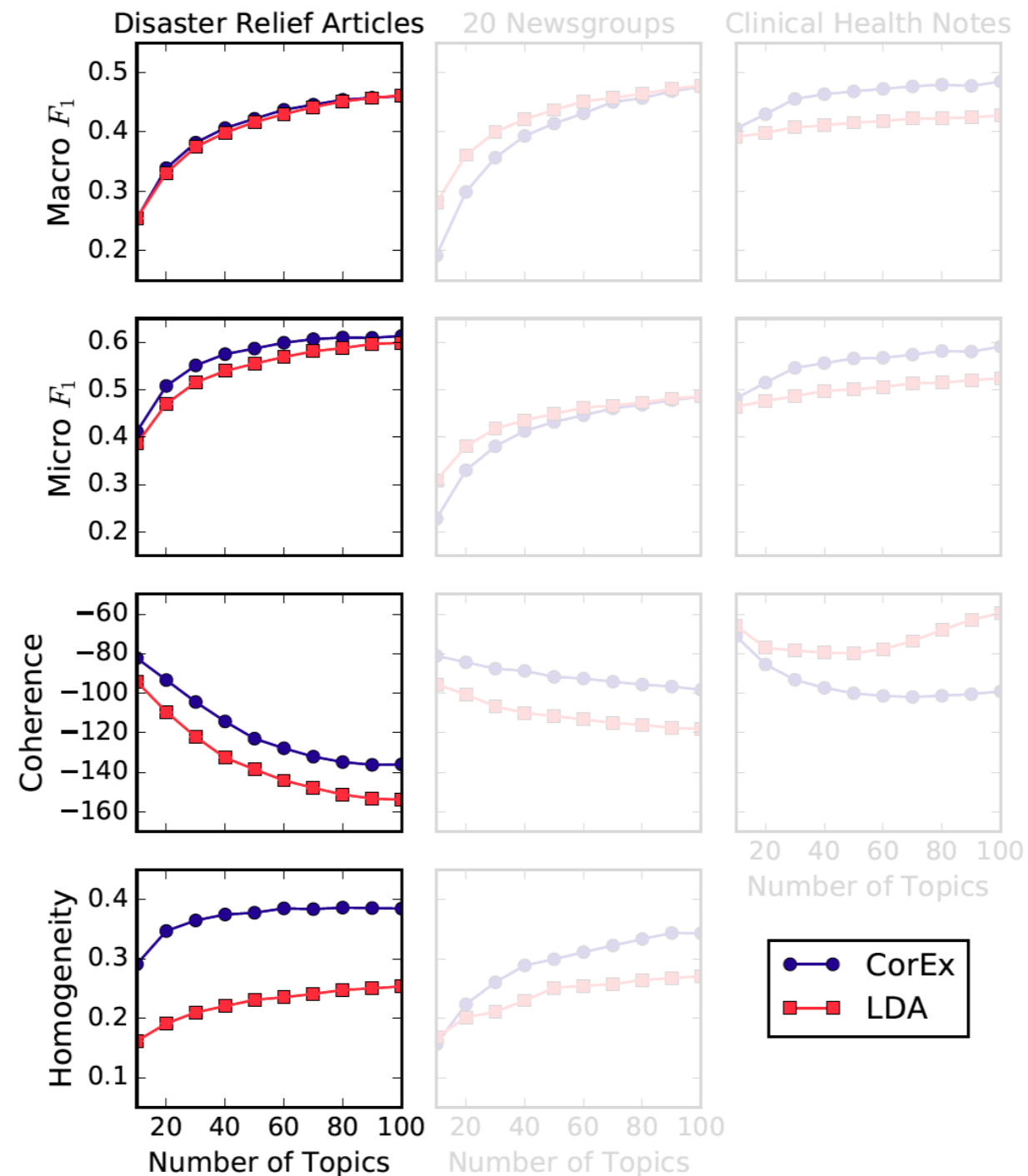
**9:** federal, its, officials, law, including, committee, staff, statement, director, group

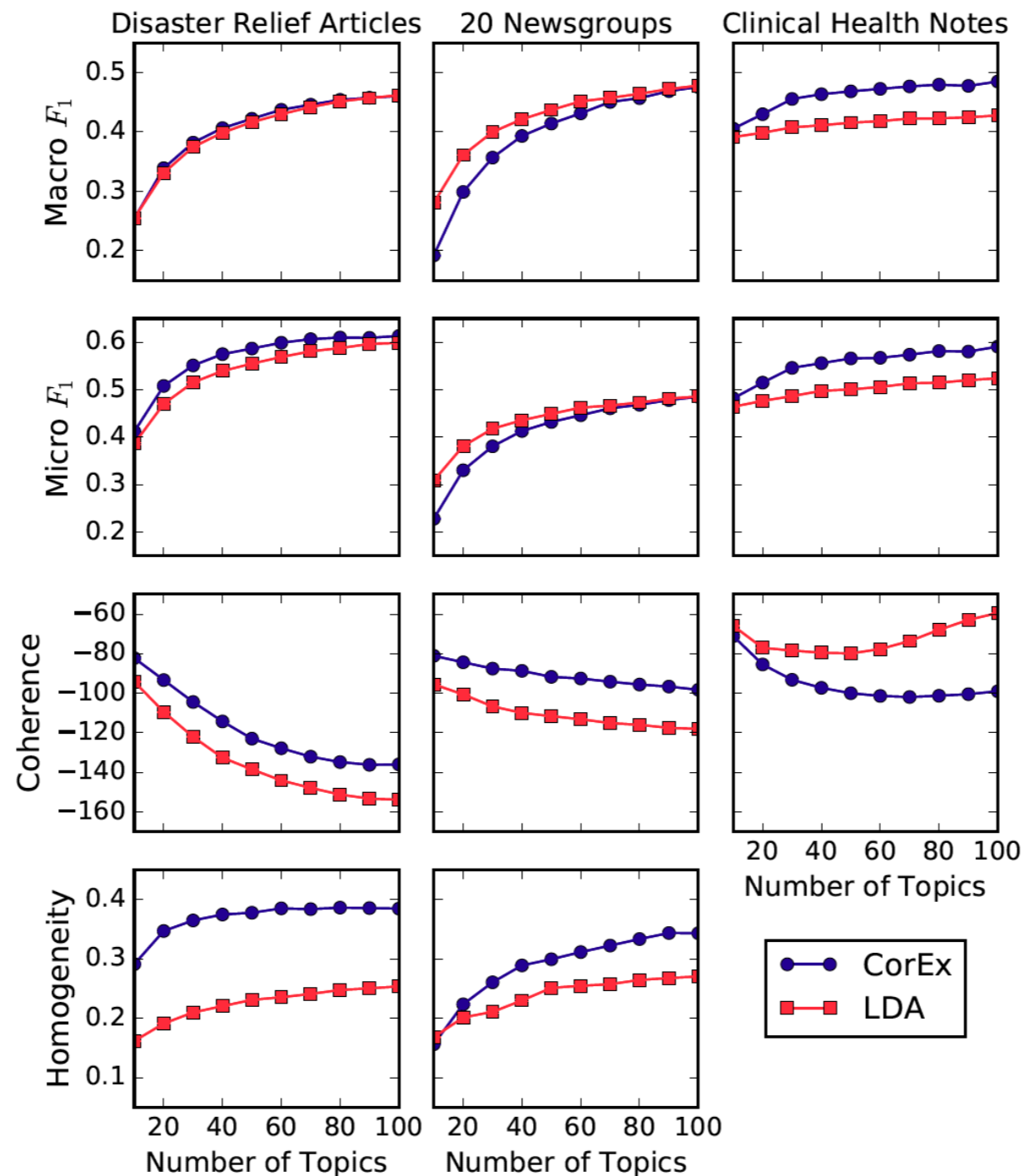**13:** islamic, foreign, military, terrorism, war, syria, iraq, isis, u, terrorist

**14:** trump, donald, trump's, republican, nominee, party, convention, top, election, him

Work by Abigail Ross and the Computational Story Lab, University of Vermont

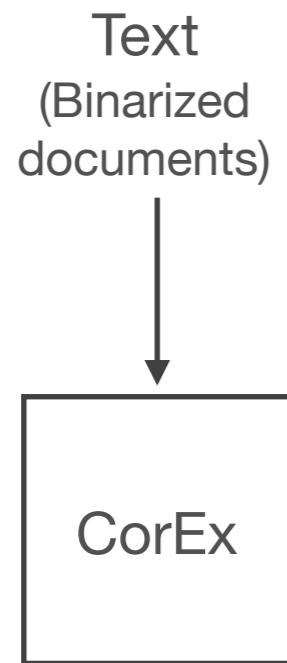# CorEx Performs Favorably Against LDA

@ryanjgallag

# CorEx Extensions

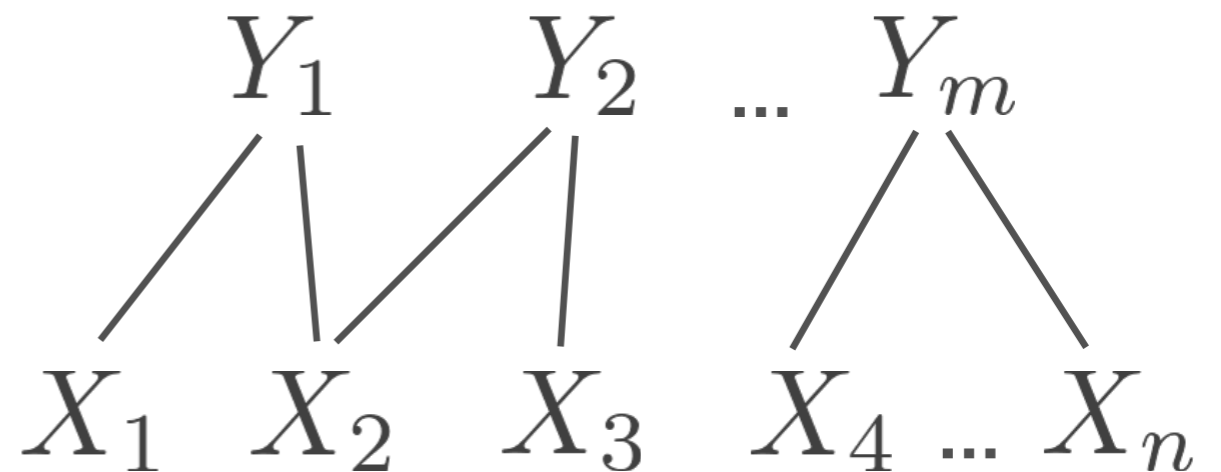With no additional assumptions, the CorEx topic model yields two extensions:
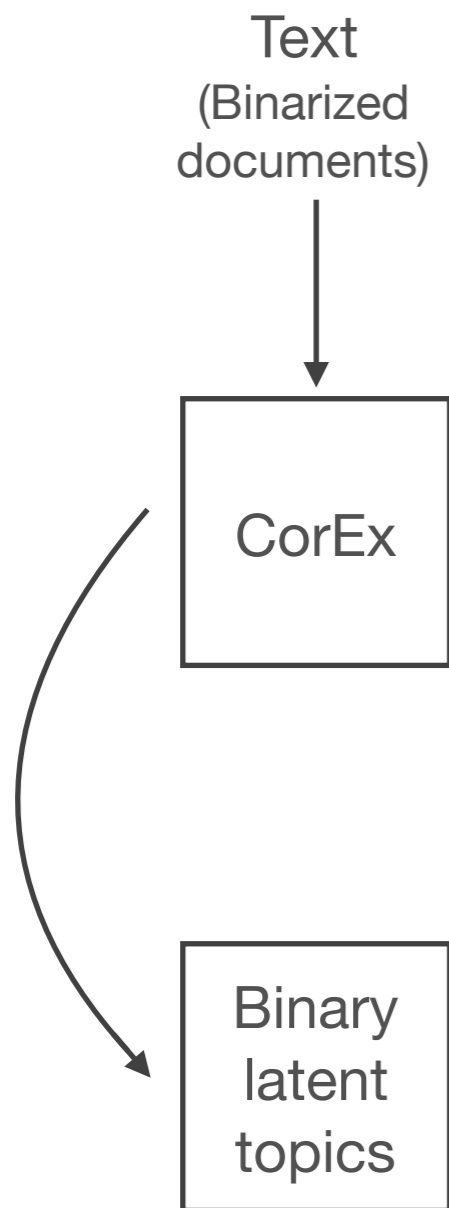
**1.** A hierarchical topic model

**2.** A semi-supervised topic model at the word level

Text
(Binarized
documents)

CorEx

$$X_1 \quad X_2 \quad X_3 \quad X_4 \ldots X_n$$

# Hierarchical CorEx

Text
(Binarized
documents)

CorEx

Binary
latent
topics

$$Y_1 \quad Y_2 \quad ... \quad Y_m$$

$$X_1 \quad X_2 \quad X_3 \quad X_4 \, ... \, X_n$$

# Hierarchical CorEx

# Hierarchical CorEx

**Data:** ~20,000 humanitarian assistance and disaster relief news articles

$$Y_1 \qquad\qquad Y_2$$

$$\alpha_{i,j}$$

$$X_1 \quad X_2 \qquad X_3 \quad X_4 \quad X_5$$

Objective: $\displaystyle\max_{G_j, p(y_j | x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} \alpha_{i,j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$

@ryanjgallag

Objective:

$$\max_{G_j, p(y_j \mid x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} \alpha_{i,j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

Maintain information about individual words

$$\textcolor{red}{\alpha_{i,j}}$$

$Y_1$     $Y_2$

$X_1$   $X_2$    $X_3$   $X_4$   $X_5$

Objective:
$$\max_{G_j, p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} \alpha_{i,j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

Maintain information about individual words      Compress documents into topics

# Anchored CorEx and the Information Bottleneck



$$Y_1 \qquad Y_2$$

$$\textcolor{red}{\alpha_{i,j}}$$

$$X_1 \quad X_2 \qquad X_3 \quad X_4 \quad X_5$$

Objective: $\displaystyle \max_{G_j, p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} \textcolor{red}{\alpha_{i,j} I(X_i : Y_j)} - \textcolor{blue}{I(X_{G_j} : Y_j)}$

Maintain information about individual words      Compress documents into topics

Information bottleneck

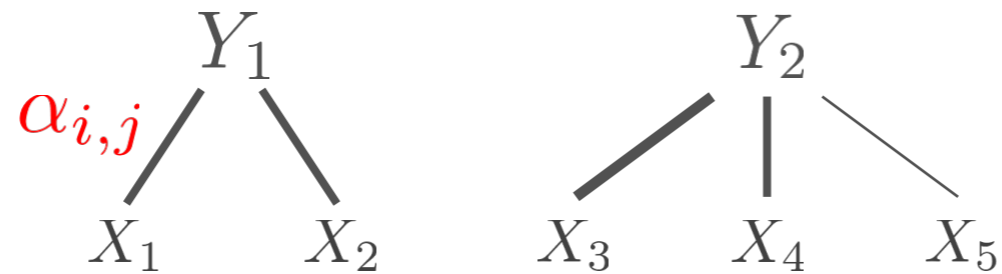"The Information Bottleneck Method." Tishby et al. (2000).

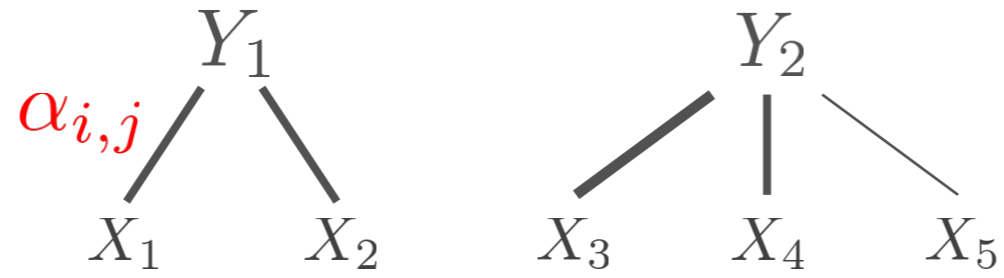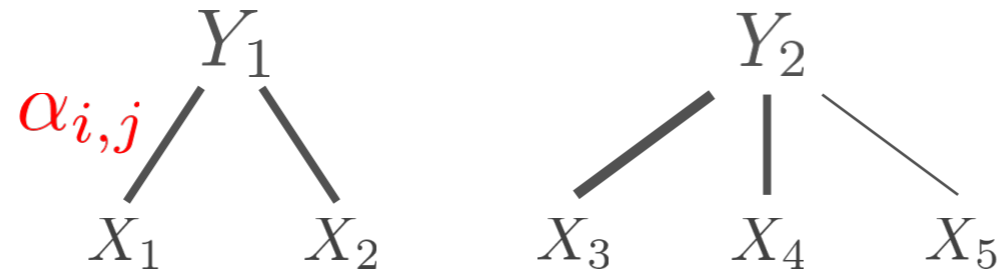# Anchored CorEx and the Information Bottleneck



Objective:
$$\max_{G_j, p(y_j|x_{G_j})} \sum_{j=1}^{m} \sum_{i \in G_j} \alpha_{i,j} I(X_i : Y_j) - I(X_{G_j} : Y_j)$$

Maintain information about individual words

Compress documents into topics

Information bottleneck

A user can **anchor** words to the latent topics by modifying the **weight** of the relationship between a word and a topic

"The Information Bottleneck Method." Tishby et al. (2000).

@ryanjgallag

# Anchoring Strategies

**Topic Representation**

Anchoring to unveil topics that do
not naturally emerge

$Y_1$

volcano    lava

$Y_2$

avalanche    snow    freezing

# Anchoring Strategies

**Topic Representation**

Anchoring to unveil topics that do not naturally emerge

$Y_1$

volcano   lava

$Y_2$

avalanche   snow   freezing

**Topic Separability**

Anchoring to help enforce separation between topics

$Y_1$

$Y_2$

computational   science   social   media   platform

# Anchoring Strategies

**Topic Representation**

Anchoring to unveil topics that do not naturally emerge

$Y_1$ — volcano, lava

$Y_2$ — avalanche, snow, freezing

**Topic Separability**

Anchoring to help enforce separation between topics

$Y_1$ — computational, science, social

$Y_2$ — social, media, platform

**Topic Aspects**

Anchoring to disambiguate different frames around a word

$Y_1$ $Y_2$ $Y_3$ $Y_4$ — election

# Anchoring for Topic Representation

**Data:** news articles about the campaigns of Clinton and Trump, up to August 2016

**Method:** train one CorEx topic model for each corpus, anchor words for comparison

$Y_1$
immigrants   immigration

$Y_2$
muslims   islam

$Y_3$
white   whites

Work by Abigail Ross and the Computational Story Lab, University of Vermont

# Anchoring for Topic Representation

**Data:** news articles about the campaigns of Clinton and Trump, up to August 2016
**Method:** train one CorEx topic model for each corpus, anchor words for comparison

$Y_1$ — immigrants  immigration

$Y_2$ — muslims  islam

$Y_3$ — white  whites

**Clinton Topic**

**1: immigration, immigrants**, jobs, economic, trade, health, tax, wall, care, economy

**Trump Topic**

**1: immigration, immigrants**, illegal, border, mexican, undocumented, rapists, mexico, wall, illegally

Work by Abigail Ross and the Computational Story Lab, University of Vermont

# Anchoring for Topic Representation

**Data:** news articles about the campaigns of Clinton and Trump, up to August 2016
**Method:** train one CorEx topic model for each corpus, anchor words for comparison
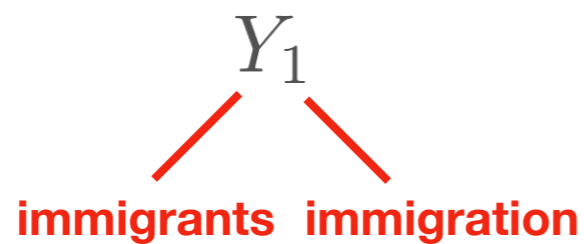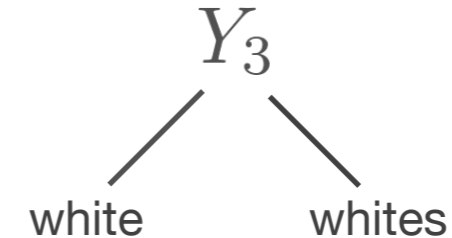
$$Y_1 \qquad\qquad Y_2 \qquad\qquad Y_3$$

immigrants    immigration        **muslims**    **islam**        white    whites

**Clinton Topic**

**2: muslims, islam**, islamic, gun, terrorism, war, military, iraq, terrorist, syria

**Trump Topic**

**2: muslims, islam**, united, ban, entering, islamic, muslim, terrorism, terrorist, terrorists

Work by Abigail Ross and the Computational Story Lab, University of Vermont

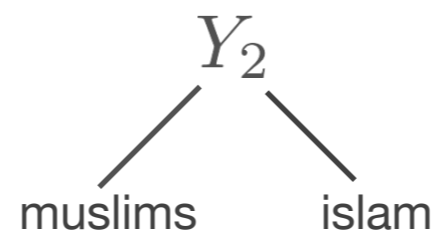@ryanjgallag

# Anchoring for Topic Representation

**Data:** news articles about the campaigns of Clinton and Trump, up to August 2016
**Method:** train one CorEx topic model for each corpus, anchor words for comparison

$$Y_1$$

immigrants    immigration

$$Y_2$$

muslims    islam

$$Y_3$$

**white**    **whites**

**Clinton Topic**

3: **white**, i, you, what, do, if, we, it's, like, people

**Trump Topic**

3: **white**, house, **whites**, supremacists, supremacist, duke, klan, klux, ku, supremacy

NAACL 2018, New Orleans, LA    @ryanjgallag

# Anchoring for Topic Aspects

**Data:** ~1 million English newswire articles since June 2015 from countries in the Middle East

$$Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \quad Y_6 \quad Y_7 \quad Y_8 \quad Y_9 \quad Y_{10}$$

aleppo

Work by Brendan Kennedy and Greg Ver Steeg, Information Sciences Institute

@ryanjgallag

**Data:** ~1 million English newswire articles since June 2015 from countries in the Middle East

$$Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \quad Y_6 \quad Y_7 \quad Y_8 \quad Y_9 \quad Y_{10}$$

aleppo

**Note:** this data broadly covers the Middle East and a priori we do not expect 10 topics to emerge about Aleppo

Work by Brendan Kennedy and Greg Ver Steeg, Information Sciences Institute

# Anchoring for Topic Aspects

**Data:** ~1 million English newswire articles since June 2015 from countries in the Middle East

$$Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \quad Y_6 \quad Y_7 \quad Y_8 \quad Y_9 \quad Y_{10}$$

aleppo

**1: aleppo**, killed, police, security, attack, state, arrested, authorities

**2: aleppo**, forces, syria, military, war, army, civilians, iraq, militants

**3: aleppo**, health, medical, food, care, water, small, conditions, treatment, patients

**4:** country, **aleppo**, east, across, group, region, middle

**5:** two, **aleppo**, took, another, place, taking, leaders

**6: aleppo**, russia, iran, barack, obama, moscow, washington, putin

**7: aleppo**, political, court, part, accused, opposition, called, saying, parliament, democratic

**8:** government, **aleppo**, minister, foreign, states, united, prime, UN, law, nations

**9: aleppo**, city, area, near, air, northern, least, town, eastern, injured

**10: aleppo**, people, children, human, rights, women, social, school, society, lives

Work by Brendan Kennedy and Greg Ver Steeg, Information Sciences Institute

NAACL 2018, New Orleans, LA                    @ryanjgallag

**Data:** ~1 million English newswire articles since June 2015 from countries in the Middle East

$$Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \quad Y_6 \quad Y_7 \quad Y_8 \quad Y_9 \quad Y_{10}$$

aleppo

**1: aleppo**, killed, police, security, attack, state, arrested, authorities

**2: aleppo**, forces, syria, military, war, army, civilians, iraq, militants

3: **aleppo**, health, medical, food, care, water, small, conditions, treatment, patients

4: country, **aleppo**, east, across, group, region, middle

5: two, **aleppo**, took, another, place, taking, leaders

6: **aleppo**, russia, iran, barack, obama, moscow, washington, putin

7: **aleppo**, political, court, part, accused, opposition, called, saying, parliament, democratic

8: government, **aleppo**, minister, foreign, states, united, prime, UN, law, nations

9: **aleppo**, city, area, near, air, northern, least, town, eastern, injured

10: **aleppo**, people, children, human, rights, women, social, school, society, lives

Work by Brendan Kennedy and Greg Ver Steeg, Information Sciences Institute

# Anchoring for Topic Aspects

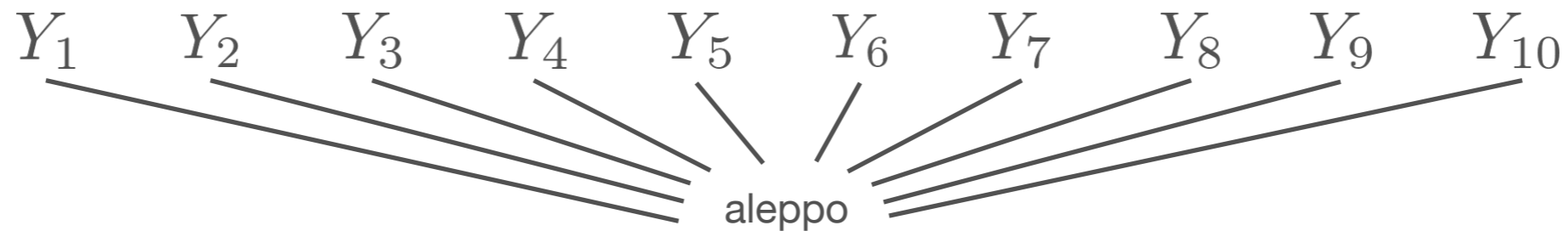**Data:** ~1 million English newswire articles since June 2015 from countries in the Middle East

$$Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \quad Y_6 \quad Y_7 \quad Y_8 \quad Y_9 \quad Y_{10}$$

aleppo

1: **aleppo**, killed, police, security, attack, state, arrested, authorities

2: **aleppo**, forces, syria, military, war, army, civilians, iraq, militants

**3: aleppo**, health, medical, food, care, water, small, conditions, treatment, patients

4: country, **aleppo**, east, across, group, region, middle

5: two, **aleppo**, took, another, place, taking, leaders

6: **aleppo**, russia, iran, barack, obama, moscow, washington, putin

7: **aleppo**, political, court, part, accused, opposition, called, saying, parliament, democratic

8: government, **aleppo**, minister, foreign, states, united, prime, UN, law, nations
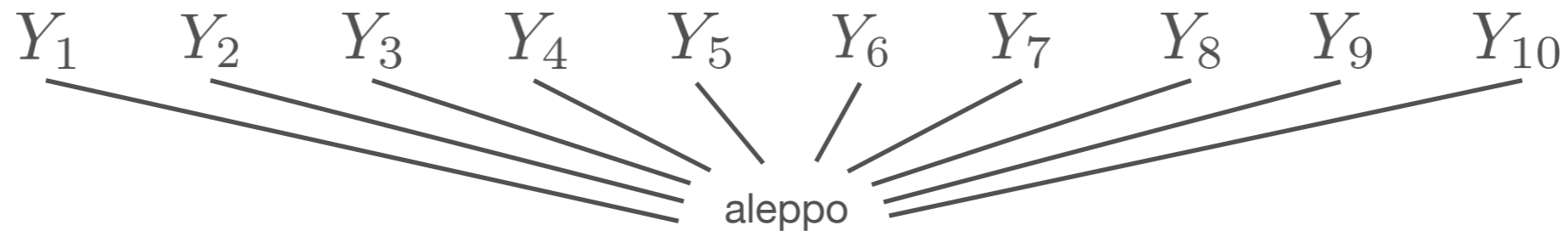
9: **aleppo**, city, area, near, air, northern, least, town, eastern, injured

10: **aleppo**, people, children, human, rights, women, social, school, society, lives

Work by Brendan Kennedy and Greg Ver Steeg, Information Sciences Institute

NAACL 2018, New Orleans, LA                    @ryanjgallag

# Anchoring for Topic Aspects

**Data:** ~1 million English newswire articles since June 2015 from countries in the Middle East

$$Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \quad Y_6 \quad Y_7 \quad Y_8 \quad Y_9 \quad Y_{10}$$

aleppo

1: **aleppo**, killed, police, security, attack, state, arrested, authorities

2: **aleppo**, forces, syria, military, war, army, civilians, iraq, militants

3: **aleppo**, health, medical, food, care, water, small, conditions, treatment, patients

4: country, **aleppo**, east, across, group, region, middle

5: two, **aleppo**, took, another, place, taking, leaders

6: **aleppo**, russia, iran, barack, obama, moscow, washington, putin

7: **aleppo**, political, court, part, accused, opposition, called, saying, parliament, democratic

8: government, **aleppo**, minister, foreign, states, united, prime, UN, law, nations
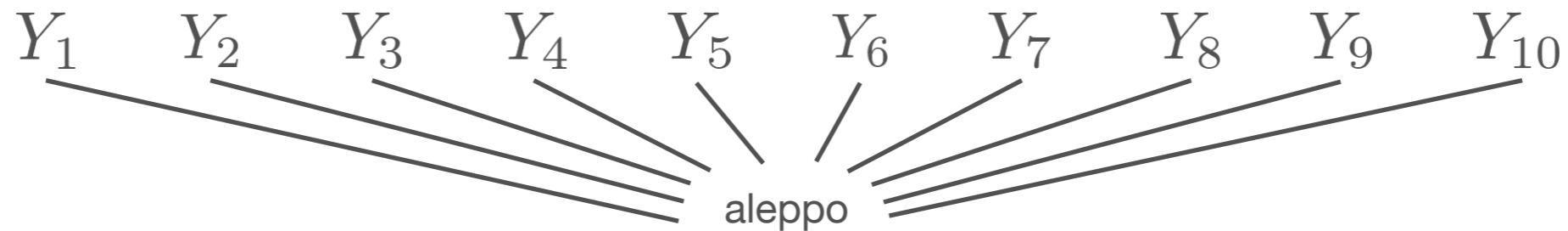
9: **aleppo**, city, area, near, air, northern, least, town, eastern, injured

10: **aleppo**, people, children, human, rights, women, social, school, society, lives

Work by Brendan Kennedy and Greg Ver Steeg, Information Sciences Institute

# Anchoring for Topic Aspects

**Data:** ~1 million English newswire articles since June 2015 from countries in the Middle East

$$Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \quad Y_6 \quad Y_7 \quad Y_8 \quad Y_9 \quad Y_{10}$$

aleppo

1: **aleppo**, killed, police, security, attack, state, arrested, authorities

2: **aleppo**, forces, syria, military, war, army, civilians, iraq, militants

3: **aleppo**, health, medical, food, care, water, small, conditions, treatment, patients

4: country, **aleppo**, east, across, group, region, middle

5: two, **aleppo**, took, another, place, taking, leaders

6: **aleppo**, russia, iran, barack, obama, moscow, washington, putin

7: **aleppo**, political, court, part, accused, opposition, called, saying, parliament, democratic

8: government, **aleppo**, minister, foreign, states, united, prime, UN, law, nations
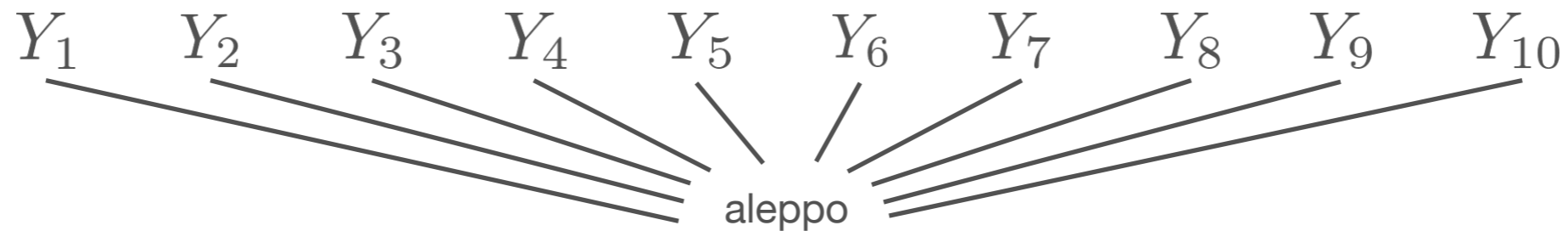
9: **aleppo**, city, area, near, air, northern, least, town, eastern, injured

10: **aleppo**, people, children, human, rights, women, social, school, society, lives

Work by Brendan Kennedy and Greg Ver Steeg, Information Sciences Institute

# Anchoring for Topic Aspects

**Data:** ~1 million English newswire articles since June 2015 from countries in the Middle East

$$Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \quad Y_6 \quad Y_7 \quad Y_8 \quad Y_9 \quad Y_{10}$$

aleppo

**1: aleppo**, killed, police, security, attack, state, arrested, authorities

**2: aleppo**, forces, syria, military, war, army, civilians, iraq, militants

**3: aleppo**, health, medical, food, care, water, small, conditions, treatment, patients

**4:** country, **aleppo**, east, across, group, region, middle

**5:** two, **aleppo**, took, another, place, taking, leaders

**6: aleppo**, russia, iran, barack, obama, moscow, washington, putin

**7: aleppo**, political, court, part, accused, opposition, called, saying, parliament, democratic

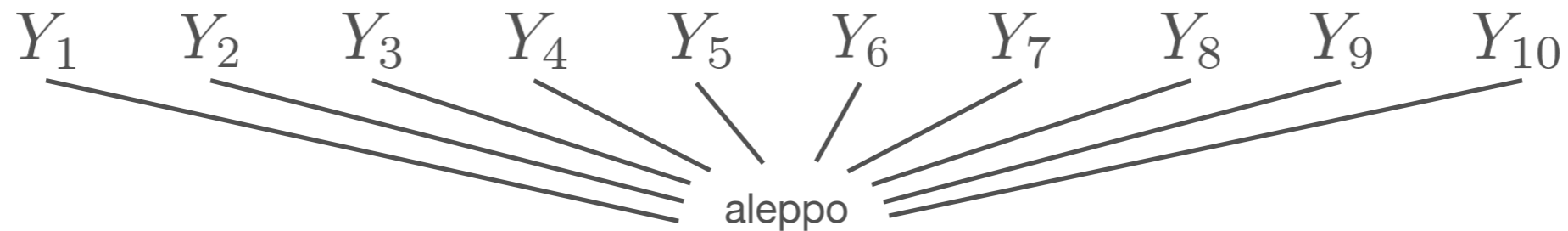**8:** government, **aleppo**, minister, foreign, states, united, prime, UN, law, nations

**9: aleppo**, city, area, near, air, northern, least, town, eastern, injured

**10: aleppo**, people, children, human, rights, women, social, school, society, lives

Work by Brendan Kennedy and Greg Ver Steeg, Information Sciences Institute

NAACL 2018, New Orleans, LA

@ryanjgallag

66

# Shape of the CorEx Topic Model to Come

### CorEx Topic Model

By defining topics in terms of information content, the CorEx topic model takes a new perspective on topic modeling

CorEx is competitive with unsupervised and semi-supervised variants of LDA while making far fewer assumptions

Anchoring through the information bottleneck provides a flexible mechanism to retrieve topics of interest and inject expert domain knowledge

### Future Work

Extend CorEx to efficiently learn multi-membership topics (*in progress*)

Incorporate count data into the CorEx topic model while preserving the benefits of the sparsity optimization

**Code:** github.com/gregversteeg/corex_topic

# Shape of the CorEx Topic Model to Come

### CorEx Topic Model

By defining topics in terms of information content, the CorEx topic model takes a new perspective on topic modeling

CorEx is competitive with unsupervised and semi-supervised variants of LDA while making far fewer assumptions

Anchoring through the information bottleneck provides a flexible mechanism to retrieve topics of interest and inject expert domain knowledge

### Future Work

Extend CorEx to efficiently learn multi-membership topics (*in progress*)

Incorporate count data into the CorEx topic model while preserving the benefits of the sparsity optimization

**Code:** github.com/gregversteeg/corex_topic

# Shape of the CorEx Topic Model to Come

## CorEx Topic Model

By defining topics in terms of information content, the CorEx topic model takes a new perspective on topic modeling

CorEx is competitive with unsupervised and semi-supervised variants of LDA while making far fewer assumptions

Anchoring through the information bottleneck provides a flexible mechanism to retrieve topics of interest and inject expert domain knowledge

## Future Work

Extend CorEx to efficiently learn multi-membership topics (*in progress*)

Incorporate count data into the CorEx topic model while preserving the benefits of the sparsity optimization

**Code:** github.com/gregversteeg/corex_topic

# Shape of the CorEx Topic Model to Come

### CorEx Topic Model

By defining topics in terms of information content, the CorEx topic model takes a new perspective on topic modeling

CorEx is competitive with unsupervised and semi-supervised variants of LDA while making far fewer assumptions

Anchoring through the information bottleneck provides a flexible mechanism to retrieve topics of interest and inject expert domain knowledge

### Future Work

Extend CorEx to efficiently learn multi-membership topics (*in progress*)

Incorporate count data into the CorEx topic model while preserving the benefits of the sparsity optimization

**Code:** github.com/gregversteeg/corex_topic

# Shape of the CorEx Topic Model to Come

### CorEx Topic Model

By defining topics in terms of information content, the CorEx topic model takes a new perspective on topic modeling

CorEx is competitive with unsupervised and semi-supervised variants of LDA while making far fewer assumptions

Anchoring through the information bottleneck provides a flexible mechanism to retrieve topics of interest and inject expert domain knowledge

### Future Work

Extend CorEx to efficiently learn multi-membership topics (*in progress*)

Incorporate count data into the CorEx topic model while preserving the benefits of the sparsity optimization

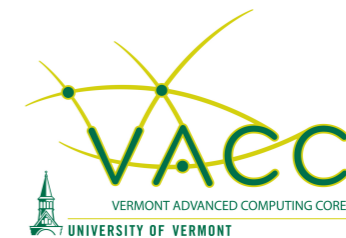**Code:** github.com/gregversteeg/corex_topic

# Collaborators

**Greg Ver Steeg**
Research Professor
Information Sciences Institute

**David Kale**
CS PhD Candidate
Information Sciences Institute

**Kyle Reing**
CS PhD Student
Information Sciences Institute

The anchored Clinton and Trump election article topics come from work by **Abigail Ross** and the **Computational Story Lab** at the University of Vermont's Complex Systems Center

# Thank you for your time!

@ryanjgallag
ryanjgallag@gmail.com

github.com/gregversteeg/corex_topic
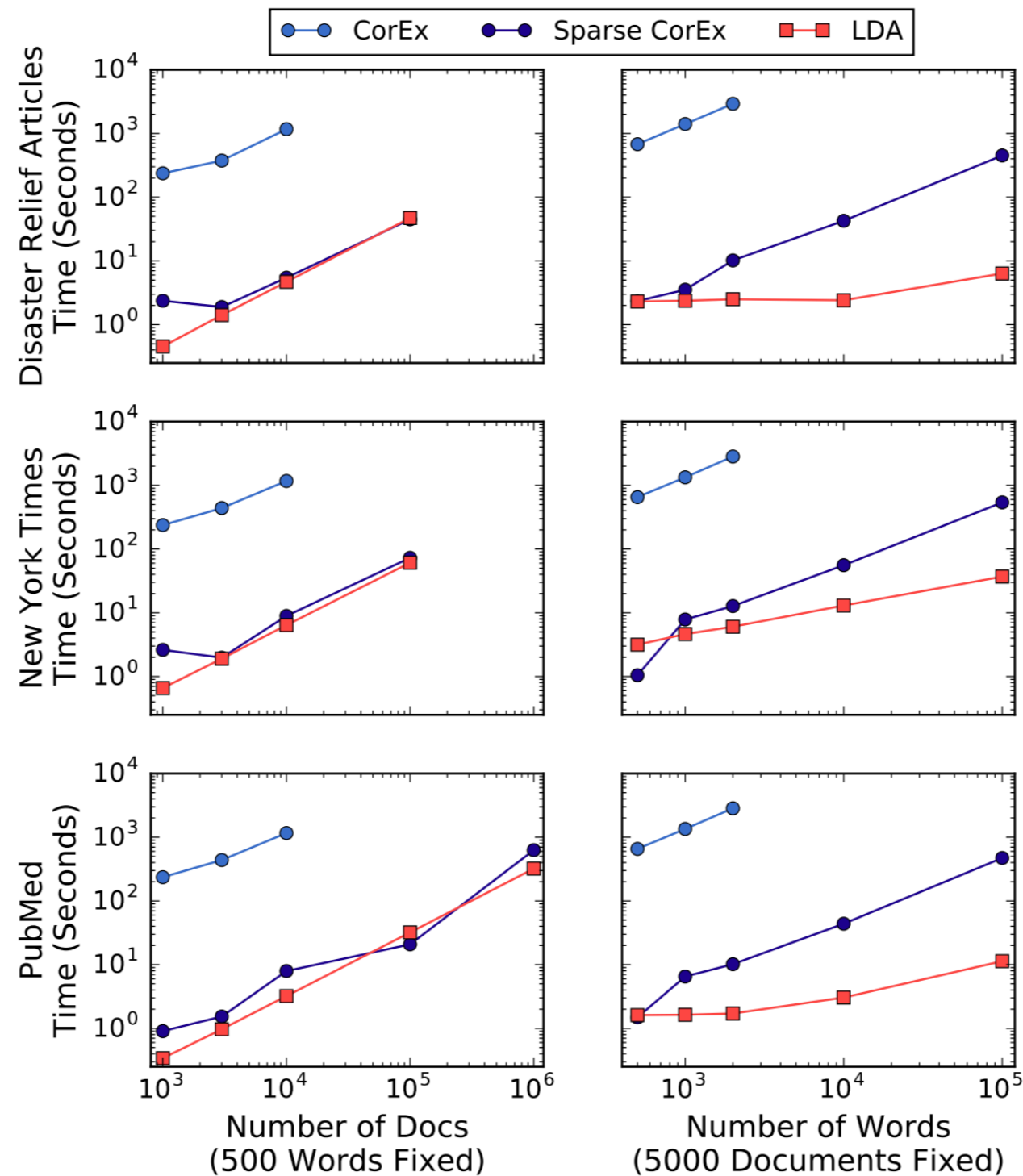
# CorEx Implementation

## Update Equations

$$p_t(y_j) = \sum_{\bar{x}} p_t(y_j \mid \bar{x}) p(\bar{x})$$

$$p_t(x_i \mid y_j) = \sum_{\bar{x}} \frac{p_t(y_j \mid \bar{x}) p(\bar{x}) \mathbb{I}[\bar{x}_i = x_i]}{p_t(y_j)}$$

Marginals in terms of the optimization parameter $p_t(y_j \mid x)$

$$\log p_{t+1}(y_j \mid x^\ell) = \log p_t(y_j) + \sum_{i=1}^{n} \alpha_{i,j}^t \log \frac{p_t(x_i^\ell \mid y_j)}{p(x_i^\ell)} - \log \mathcal{Z}_j(x^\ell)$$

Probabilistic labels for each latent factor given sample

## Sparsity Optimization

$$\log \frac{p_t(x^\ell \mid y_j)}{p(x_i^\ell)} = \log \frac{p_t(X_i = 0 \mid y_j)}{p(X_i = 0)} + x_i^\ell \log \frac{p_t(X_i^\ell = 1 \mid y_j) p(X_i = 0)}{p_t(X_i = 0 \mid y_j) p(x_i^\ell = 1)}$$

Substituting above turns the sum into a matrix multiplication between a matrix of size (# docs) x (# types) and a matrix of size (# types) x (# topics)

# Sparsity Optimization Speed Comparison



CorEx ● Sparse CorEx ● LDA ■

Disaster Relief Articles Time (Seconds)

New York Times Time (Seconds)

PubMed Time (Seconds)

Number of Docs (500 Words Fixed)

Number of Words (5000 Documents Fixed)

# CorEx Example Topics

**Data:** news articles about Clinton and Trump, train one CorEx topic model for each corpus

### Clinton Article Topics

**1:** server, department, classified, information, private, investigation, fib, email, emails, secretary

**3:** sanders, bernie, primary, vermont, win, voters, race, nomination, vote, polls

**8:** percent, poll, points, percentage, margin, survey, according, 10, polling, university

**9:** federal, its, officials, law, including, committee, staff, statement, director, group

**13:** islamic, foreign, military, terrorism, war, syria, iraq, isis, u, terrorist

**14:** trump, donald, trump's, republican, nominee, party, convention, top, election, him

### Trump Article Topics

**1:** primary, party, win, cruz, delegates, voters, ted, nomination, republicans, vote

**4:** $, tax, money, million, jobs, economic, companies, billion, pay, taxes

**7:** percent, poll, percentage, points, polls, survey, 10, polling, margin, according
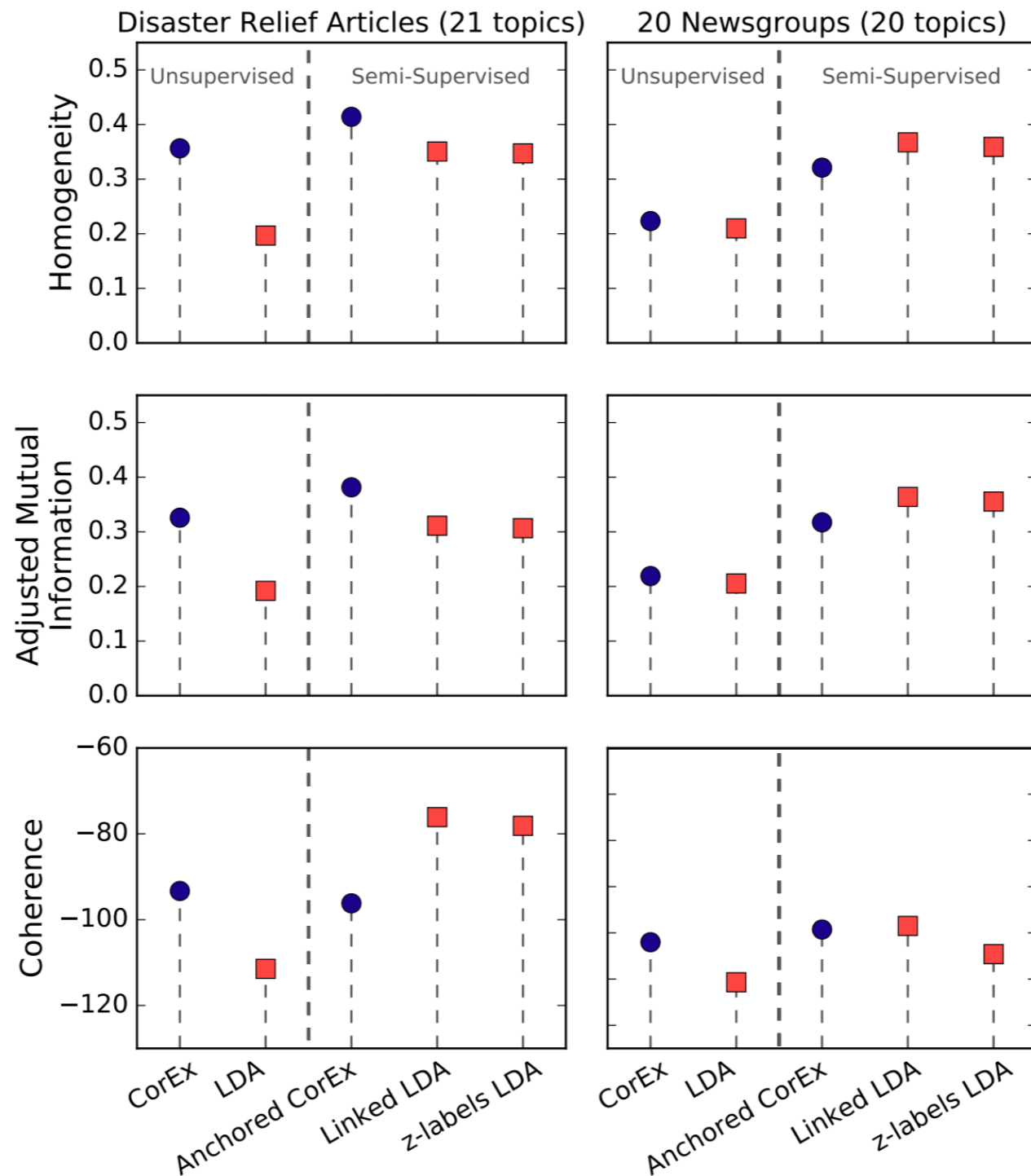
**12:** crowd, rally, night, event, speech, stage, audience, spoke, wife, took

**14:** rubio, marco, jeb, bush, carson, florida, ben, candidates, iowa, gov

**25:** clinton, hillary, bernie, sanders, democratic, clinton's, her, she, vermont, secretary

@ryanjgallag

# Comparisons to Semi-Supervised LDA



Disaster Relief Articles (21 topics)      20 Newsgroups (20 topics)

# Anchoring Experiment

**Data:** HA/DR news articles and clinical health notes

For each document label:



Determine anchor words by measuring the words with the highest mutual information with the label

Anchor one topic of CorEx topic model with the label anchor words

Run an unsupervised CorEx topic model with the same random seed

Repeat 30 times

Compute the difference in the metric as a matched pair

Analyze the distribution of the metric across models

# Anchoring Experiment: Effect of Parameter

Anchoring Experiment: Heterogeneity of Effects

NAACL 2018, New Orleans, LA

@ryanjgallag

# Anchoring for Topic Aspects

**Data:** ~870,000 unique tweets containing #Ferguson from Aug. 9th-Nov. 30th, 2014

$Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \qquad Y_6 \quad Y_7 \quad Y_8 \quad Y_9 \quad Y_{10}$

protest
protests

riot
riots

### "protest" Topics

**1: protest, protests**, peaceful, violent, continue, night, island, photos, staten, nights

**2: protest, protests**, #hiphopmoves, #cole, hiphop, nationwide, moves, fo, anheuser, boeing

**3: protest, protests**, st, louis, guard, national, county, patrol, highway, city

**4: protest, protests**, paddy, covering, beverly, walmart, wagon, hills, passionately, including

**5: protest, protests**, solidarity, march, square, rally, #oakland, downtown, nyc, #nyc

### "riot" Topics

**6: riot, riots**, unheard, language, inciting, accidentally, jokingly, watts, waving, dies

**7: riot**, black, **riots**, white, #tcot, blacks, men, whites, race, #pjnet

**8: riot, riots**, looks, like, sounds, acting, act, animals, looked, treated

**9: riot, riots**, store, looting, businesses, burning, fire, looted, stores, business

**10:** gas, **riot**, tear, **riots**, gear, rubber, bullets, military, molotov, armored

@ryanjgallag