# Spherical Admixture Models

Joseph Reisinger*, Austin Waters, Bryan Silverthorn, and Raymond J. Mooney

June 22nd, 2010

- "I want to use LDA..."

- Want to use some set of feature weights capturing semantic content (tf-idf, pmi, etc)

- Empirical benefits to cosine distance in classical IR tasks.

Dhillon and Modha (2001), Strehl et al. (2000), Salton and McGill (1983)

$$
\begin{array}{llll}
\boldsymbol{\theta}_d|\boldsymbol{\alpha} & \sim & \text{Dirichlet}(\boldsymbol{\alpha}), & d \in D, & \text{(topic proportions)} \\
\boldsymbol{\phi}_t|\boldsymbol{\beta} & \sim & \text{Dirichlet}(\boldsymbol{\beta}), & t \in T, & \text{(topics)} \\
z_{id}|\boldsymbol{\theta}_d & \sim & \text{Mult}(\boldsymbol{\theta}_d), & i \in |\mathbf{w}_d|, & \text{(topic indicators)} \\
w_{id}|\boldsymbol{\phi}_{z_{id}} & \sim & \text{Mult}(\boldsymbol{\phi}_{z_{id}}), & i \in |\mathbf{w}_d|, & \text{(words)}
\end{array}
$$

| $\phi_1$ | $\phi_2$ | $\phi_3$ |
|---|---|---|
| government | wrote | finance |
| minister | said | economists |
| state | responding | spending |
| federal | editor | budget |

$\sim \text{Dir}(\beta)$

$\phi =$

$d_1 =$ Responding to finance minister Ruth Richardson's May 1991 budget which cut government spending, 15 academic economists from the University of Auckland wrote a letter to the editor of the New Zealand Herald on 6 June 1991. It read: "We wish to state in the strongest possible terms our view that in the present

$\mathbf{d} =$ 「我在人生中存在的意義究竟是

$$\begin{aligned}
\boldsymbol{\theta}_d | \boldsymbol{\alpha} &\sim \mathrm{Dirichlet}(\boldsymbol{\alpha}), & d \in D, & \quad \text{(topic proportions)} \\
\boldsymbol{\phi}_t | \boldsymbol{\beta} &\sim \mathrm{Dirichlet}(\boldsymbol{\beta}), & t \in T, & \quad \text{(topics)} \\
z_{id} | \boldsymbol{\theta}_d &\sim \mathrm{Mult}(\boldsymbol{\theta}_d), & i \in |\mathbf{w}_d|, & \quad \text{(topic indicators)} \\
w_{id} | \boldsymbol{\phi}_{z_{id}} &\sim \mathrm{Mult}(\boldsymbol{\phi}_{z_{id}}), & i \in |\mathbf{w}_d|, & \quad \text{(words)}
\end{aligned}$$

| $\boldsymbol{\phi}_1$ | $\boldsymbol{\phi}_2$ | $\boldsymbol{\phi}_3$ |
|---|---|---|
| government | wrote | finance |
| minister | said | economists |
| state | responding | spending |
| federal | editor | budget |

$\sim \mathrm{Dir}(\beta)$

$\boldsymbol{\phi} =$

$d_1 =$ Responding to finance minister Ruth Richardson's May 1991 budget which cut government spending, 15 academic economists from the University of Auckland wrote a letter to the editor of the New Zealand Herald on 6 June 1991. It read: "We wish to state in the strongest possible terms our view that in the present

$\mathbf{d} =$ 「我 在人生中 存在的意 義究竟是

4

$$\begin{aligned}
\boldsymbol{\theta}_d | \boldsymbol{\alpha} &\sim \text{Dirichlet}(\boldsymbol{\alpha}), & d \in D, & \quad \text{(topic proportions)} \\
\boldsymbol{\phi}_t | \boldsymbol{\beta} &\sim \text{Dirichlet}(\boldsymbol{\beta}), & t \in T, & \quad \text{(topics)} \\
z_{id} | \boldsymbol{\theta}_d &\sim \text{Mult}(\boldsymbol{\theta}_d), & i \in |\mathbf{w}_d|, & \quad \text{(topic indicators)} \\
w_{id} | \boldsymbol{\phi}_{z_{id}} &\sim \text{Mult}(\boldsymbol{\phi}_{z_{id}}), & i \in |\mathbf{w}_d|, & \quad \text{(words)}
\end{aligned}$$

- Topic modeling is basically the same story as dimensionality reduction, e.g. SVD, PCA, NMF, ...

- Differences:

  - Bayesian

  - More emphasis on interpreting topics

  - Generative models offer more flexibility

$$
\begin{aligned}
\boldsymbol{\theta}_d | \boldsymbol{\alpha} &\sim \text{Dirichlet}(\boldsymbol{\alpha}), & d \in D, && \text{(topic proportions)} \\
\boldsymbol{\phi}_t | \boldsymbol{\beta} &\sim \text{Dirichlet}(\boldsymbol{\beta}), & t \in T, && \text{(topics)} \\
z_{id} | \boldsymbol{\theta}_d &\sim \text{Mult}(\boldsymbol{\theta}_d), & i \in |\mathbf{w}_d|, && \text{(topic indicators)} \\
w_{id} | \boldsymbol{\phi}_{z_{id}} &\sim \text{Mult}(\boldsymbol{\phi}_{z_{id}}), & i \in |\mathbf{w}_d|, && \text{(words)}
\end{aligned}
$$

- We can explicitly represent the multinomial distribution that a document is drawn from integrating out z instead of theta:

$$
w_{id} \sim \text{Mult}(\theta_d^\top \Phi)
$$

- i.e. a weighted average over the topics.

(Blei et al. 2003)

# Spherical mixture modeling intuition

spherical mixture model

$$\phi_k \sim \mathrm{vMF}(\mathbf{m}_0) \qquad k \in K \qquad \text{(clusters)}$$
$$z_i \sim H \qquad\qquad i \in D \qquad \text{(assignments)}$$
$$\boldsymbol{d}_i \sim \mathrm{vMF}(\phi_{z_i}) \qquad i \in D \qquad \text{(documents)}$$

von Mises-Fisher Distribution
$$f(\mathbf{x}; \boldsymbol{\mu}, \kappa) = c_d(\kappa) \exp\left(\kappa \boldsymbol{\mu}^\top \mathbf{x}\right)$$
$$||\boldsymbol{\mu}|| = 1,\ \kappa \geq 0$$

- Generalization of spherical k-means / cosine distance

- Embed documents in the unit-hypersphere (L2 norm)

- Cosine distance has been quite successful in IR / document modeling (less sensitive to any one single feature)

(Banerjee et al. 2006)

$$\begin{aligned}
\boldsymbol{\theta}_d | \boldsymbol{\alpha} &\sim \text{Dirichlet}(\boldsymbol{\alpha}), & d \in D, & \quad \text{(topic proportions)} \\
\boldsymbol{\phi}_t | \boldsymbol{\beta} &\sim \text{Dirichlet}(\boldsymbol{\beta}), & t \in T, & \quad \text{(topics)} \\
z_{id} | \boldsymbol{\theta}_d &\sim \text{Mult}(\boldsymbol{\theta}_d), & i \in |\mathbf{w}_d|, & \quad \text{(topic indicators)} \\
w_{id} | \boldsymbol{\phi}_{z_{id}} &\sim \text{Mult}(\boldsymbol{\phi}_{z_{id}}), & i \in |\mathbf{w}_d|, & \quad \text{(words)}
\end{aligned}$$

$+$

## Spherical mixture model

$$\begin{aligned}
\boldsymbol{\phi}_k &\sim \text{vMF}(\boldsymbol{\mu}, \xi) & k \in K & \quad \text{(clusters)} \\
z_i &\sim H & i \in D & \quad \text{(assignments)} \\
\boldsymbol{d}_i &\sim \text{vMF}(\boldsymbol{\phi}_{z_i}) & i \in D & \quad \text{(documents)}
\end{aligned}$$

$=$

## Spherical Admixture Model

$$\begin{aligned}
\boldsymbol{\mu} | \kappa_0 &\sim \text{vMF}(\mathbf{m}, \kappa_0), & & \text{(corpus mean)} \\
\boldsymbol{\phi}_t | \boldsymbol{\mu}, \xi &\sim \text{vMF}(\boldsymbol{\mu}, \xi), & t \in T, & \text{(topics)} \\
\boldsymbol{\theta}_d | \boldsymbol{\alpha} &\sim \text{Dirichlet}(\boldsymbol{\alpha}), & d \in D, & \text{(topic proportions)} \\
\bar{\boldsymbol{\phi}}_d | \boldsymbol{\phi}, \boldsymbol{\theta}_d &= \text{Avg}(\boldsymbol{\phi}, \boldsymbol{\theta}_d), & d \in D, & \text{(spherical average)} \\
\mathbf{v}_d | \bar{\boldsymbol{\phi}}_d, \kappa &\sim \text{vMF}\left(\bar{\boldsymbol{\phi}}_d, \kappa\right), & d \in D, & \text{(documents)}
\end{aligned}$$

$$\begin{aligned}
\boldsymbol{\mu}|\kappa_0 &\sim \\
\phi_t|\boldsymbol{\mu},\xi &\sim \\
\boldsymbol{\theta}_d|\boldsymbol{\alpha} &\sim \\
\bar{\phi}_d|\phi,\boldsymbol{\theta}_d &= \\
\mathbf{v}_d|\bar{\phi}_d,\kappa &\sim
\end{aligned}$$

- LDA

- This
  whe

- So w
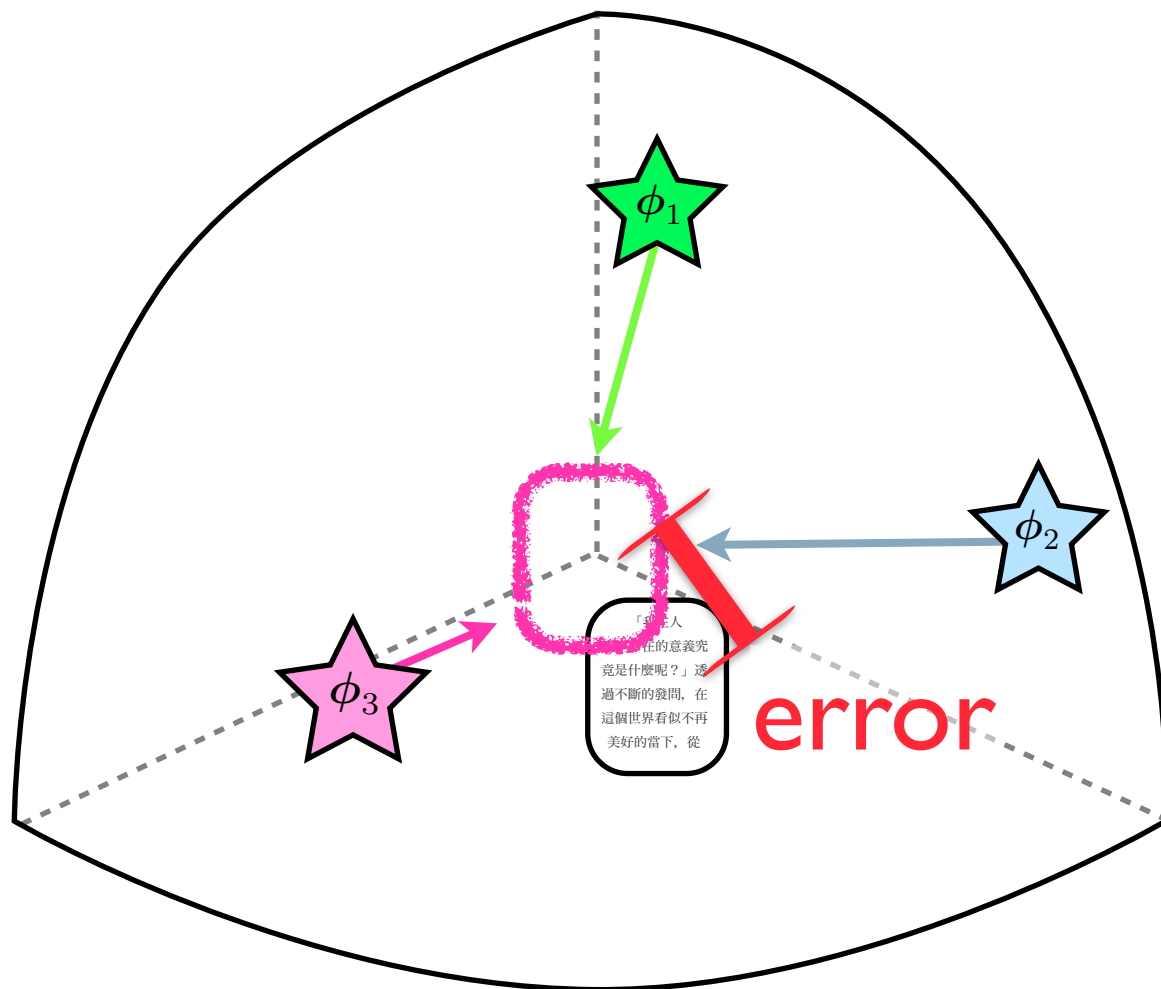
# Drawing documents



Latent Dirichlet Allocation

Spherical Admixture Model

# Drawing documents

## Spherical Admixture Model

$$
\begin{array}{llll}
\boldsymbol{\mu}|\kappa_0 & \sim & \text{vMF}(\mathbf{m},\kappa_0), & & \text{(corpus mean)} \\
\phi_t|\boldsymbol{\mu},\xi & \sim & \text{vMF}(\boldsymbol{\mu},\xi), & t \in T, & \text{(topics)} \\
\boldsymbol{\theta}_d|\boldsymbol{\alpha} & \sim & \text{Dirichlet}(\boldsymbol{\alpha}), & d \in D, & \text{(topic proportions)} \\
\bar{\phi}_d|\phi,\boldsymbol{\theta}_d & = & \text{Avg}(\phi,\boldsymbol{\theta}_d), & d \in D, & \text{(spherical average)} \\
\mathbf{v}_d|\bar{\phi}_d,\kappa & \sim & \text{vMF}\left(\bar{\phi}_d,\kappa\right), & d \in D, & \text{(documents)}
\end{array}
$$



$\phi_1$  $\phi_2$  $\phi_3$  error

- Variational EM for inference

- Tractable: ~10k docs in O(hours)

http://www.cs.utexas.edu/~austin

11

# Topic interpretability



**NIPS**

| (+) | (−) | (+) | (−) |
|-----|-----|-----|-----|
| svm | network | genetic | mlp |

**Wikipedia**

| (+) | (−) | (+) | (−) | (+) | (−) |
|-----|-----|-----|-----|-----|-----|
| navy | airport | album | opera | india | germany |
| ships | airlines | label | actor | temple | borough |
| naval | flights | singles | films | dynasty | england |
| submarines | bus | chart | players | indian | france |
| aircraft | satellites | song | conservatory | khan | parish |

- Observing a term with negative weight is evidence *against* that topic

- Negative weight terms are often semantically similar, near-neighbor topics

12

# Human studies: topic coherence

**LDA**

| |
|---|
| male, mammals, <u>empire</u>, plants, species, birds |
| court, crimes, police, law, security, <u>jazz</u> |

**SAM**

| |
|---|
| vishnu, tamil, kerala, singh, <u>meteorologist</u>, nadu |
| oxidation, <u>footballers</u>, protein, potassium, hydrogen, symptoms |

- Measure semantic coherence of the highest weighted terms in each topic via a "word intrusion" task

- Human raters were recruited using Mechanical Turk

- Quality control: (1) manually constructed tasks, (2) screening for low LOO inter-annotator agreement

(Chang et al. 2009)

13

# Human studies: topic coherence



- 8 raters per question (632 unique), 50 questions per model

- LDA: 52%, SAM tf 80%, SAM tf-idf 82% identification rate

(Chang et al. 2009)

# Human studies: topic relevance



DIRECTIONS: Which list best describes the main theme of the wikipedia article appearing below?

Rolls-Royce Spey

- Forced choice: "which set of words best describes the main theme of the article?"

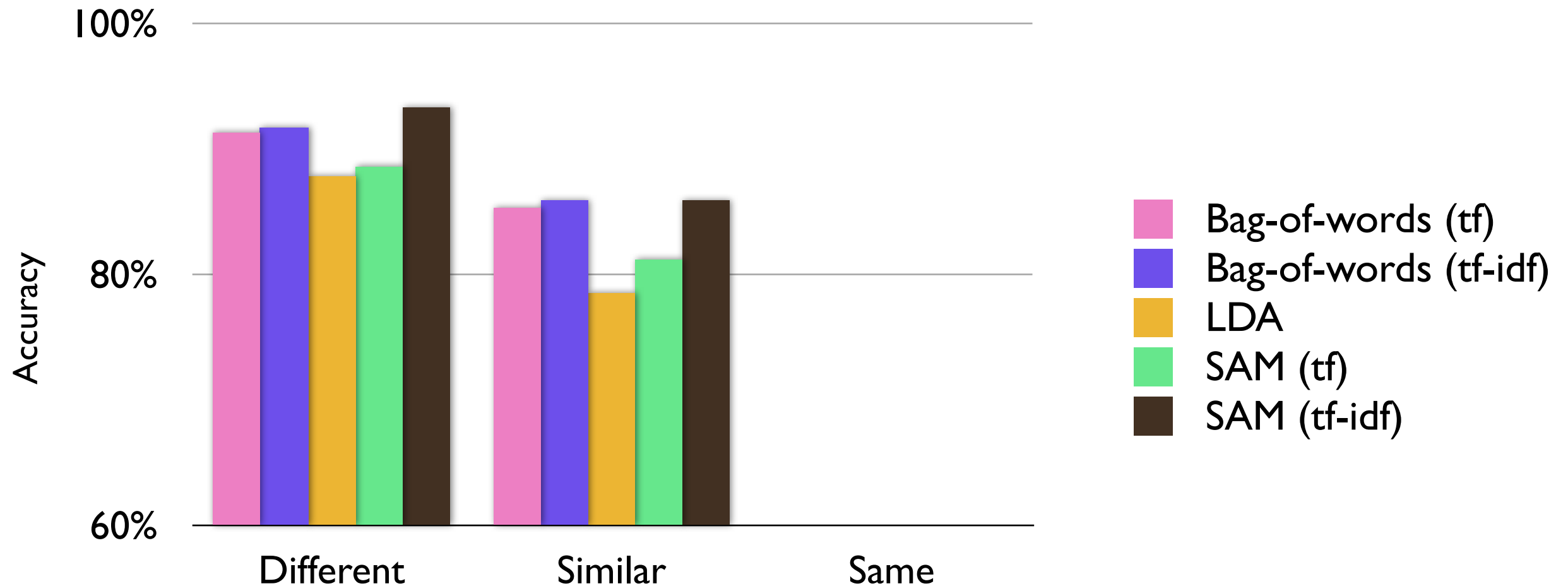- Discarded 47 articles with low kappa; SAM results preferred 62%

(Chang et al. 2009)
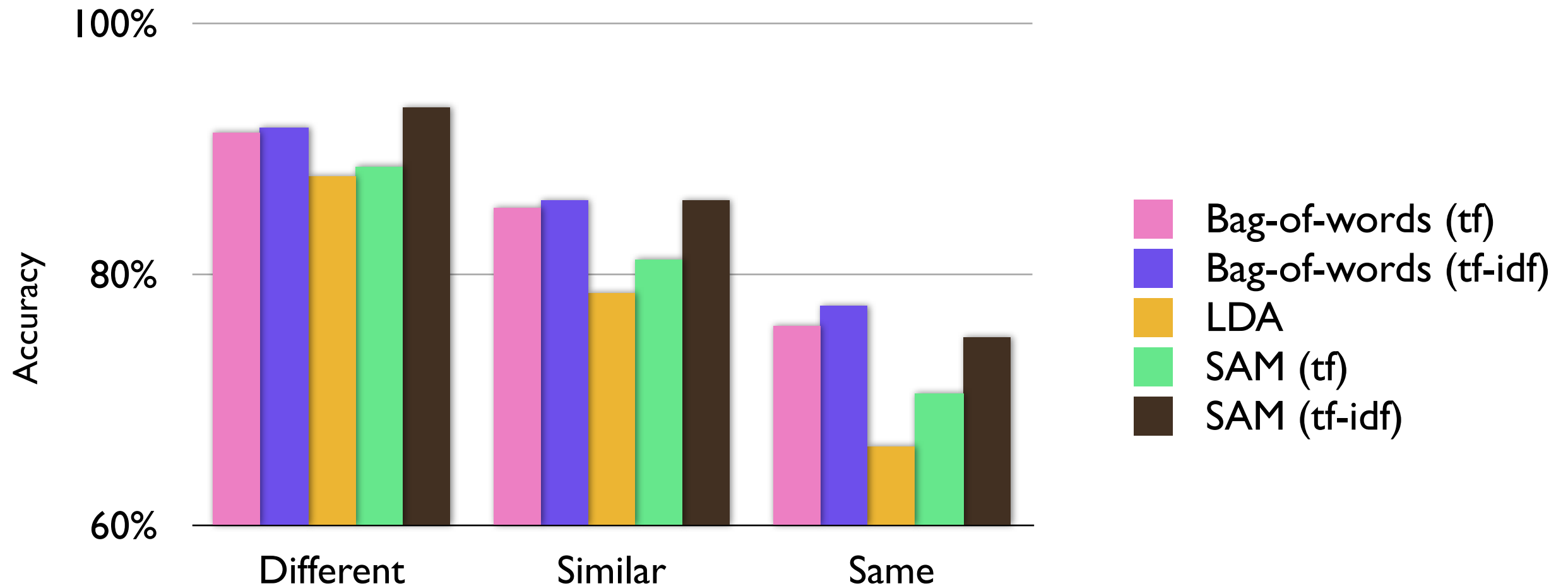
# Results: 20 newsgroups



Legend:
- Bag-of-words (tf)
- Bag-of-words (tf-idf)
- LDA
- SAM (tf)
- SAM (tf-idf)

Y-axis: Accuracy (60%, 80%, 100%)
X-axis: Different, Similar, Same

- Three classification tasks:

  - Different: `rec.sport.baseball, sci.space, alt.atheism`

  - Similar: `rec.sport.baseball, talk.politics.guns, talk.politics.misc`

  - Same: `comp.os.ms-windows.misc, comp.windows.x comp.graphics`

(Banerjee and Basu 2007)

# Results: 20 newsgroups
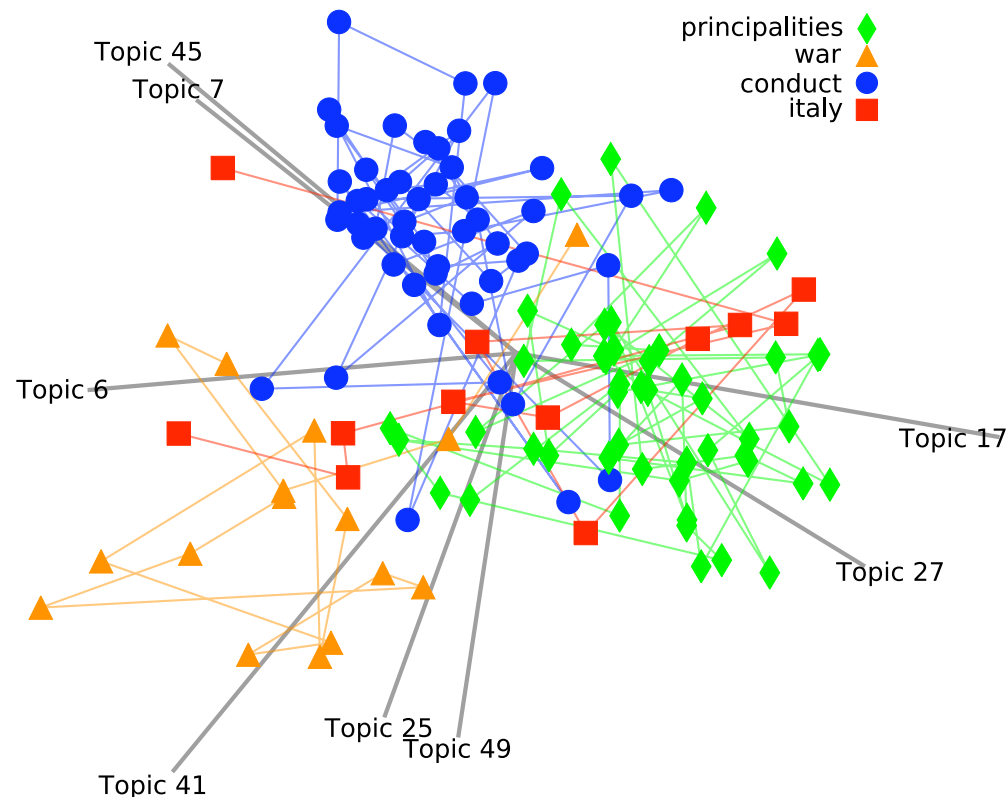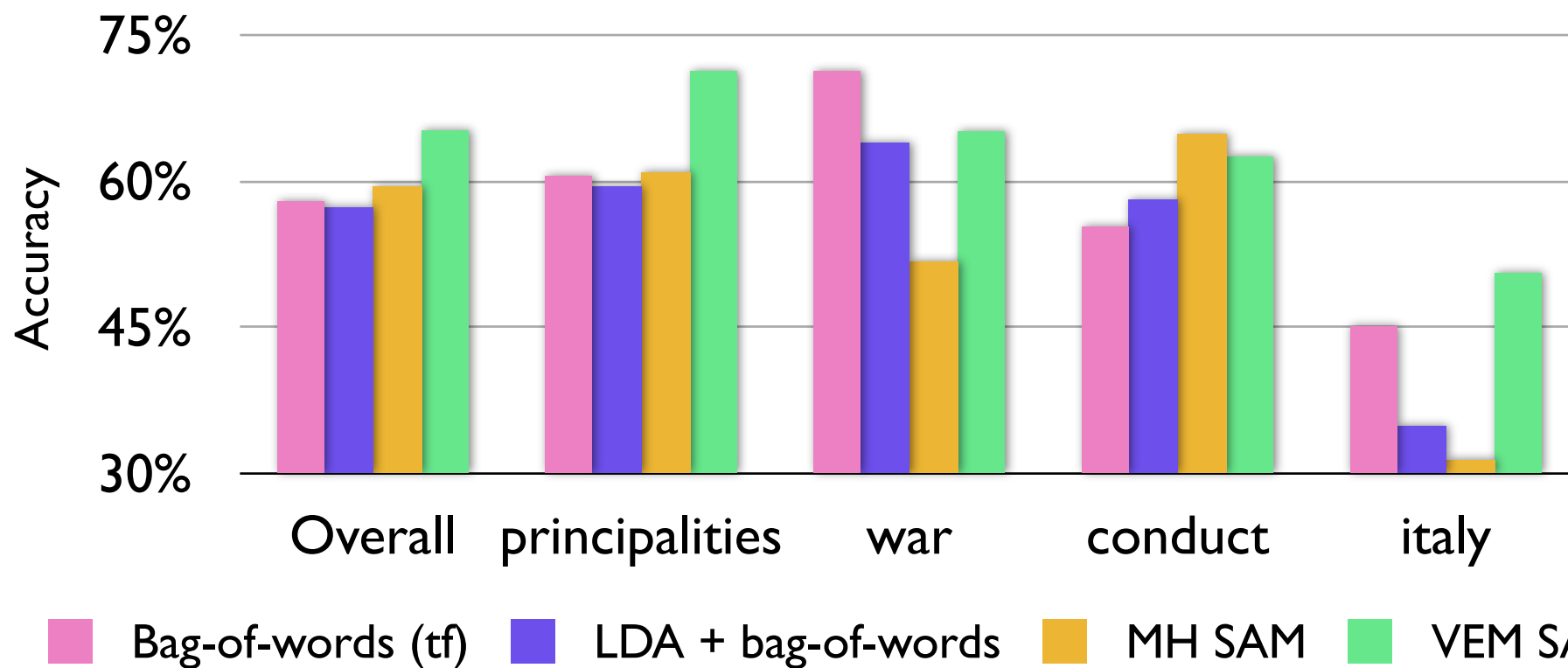


- Three classification tasks:

  - Different: `rec.sport.baseball, sci.space, alt.atheism`

  - Similar: `rec.sport.baseball, talk.politics.guns, talk.politics.misc`

  - Same: `comp.os.ms-windows.misc, comp.windows.x comp.graphics`

(Banerjee and Basu 2007)

# Results: 20 newsgroups



- Three classification tasks:

  - Different: `rec.sport.baseball, sci.space, alt.atheism`

  - Similar: `rec.sport.baseball, talk.politics.guns, talk.politics.misc`

  - Same: `comp.os.ms-windows.misc, comp.windows.x comp.graphics`

(Banerjee and Basu 2007)

# Results: *il principe*



Legend: Bag-of-words (tf) — LDA + bag-of-words — MH SAM — VEM SAM



principalities
war
conduct
italy

Topic 45
Topic 7
Topic 6
Topic 17
Topic 27
Topic 25
Topic 49
Topic 41

- Short, singly-authored, thematically tight

- 4 main themes corresponding to 4 sections:

  - Types of Principalities, Ch 1-11

  - Types of Armies, Ch 12-14

  - The Conduct of Princes, Ch 15-23

  - Political Situation in Italy, Ch 24-26

# Why does it work?

- Feature weighting helps dimensionality reduction (less for interpretability).

- Dense topic vectors can account for missing terms.

- Cosine distance may better measure document / topic similarity.

# Conclusions

- Replacing multinomial likelihood of LDA with vMF (spherical); inference is tractable

- Cosine distance; dense topic vectors

- Better results as a dimensionality reduction method

- Top weighted terms are more semantically coherent (human raters)

- Benefits are less pronounced for denser data sets (e.g. vision)

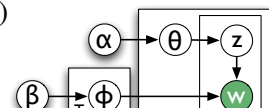- Negative weight terms capture some useful structure.

# Thanks!

$(tf)$    SAM (tf-idf)

singh, nadu, **meteorologist**
assium, **footballers**,

| | ton, county, mississippi, wl |
| (hard) | 2: tang, hong, **howe**, wu, kong, leone |
| LDA | 1: male, mammals, **empire**, plants, species, birds |
| (easy) | 2: court, crimes, police, law, security, **jazz** |

$$\theta_d | \boldsymbol{\alpha} \sim \text{Dirichlet}(\boldsymbol{\alpha}), \quad d \in D, \quad \text{(topic proportions)}$$
$$\phi_t | \boldsymbol{\beta} \sim \text{Dirichlet}(\boldsymbol{\beta}), \quad t \in T, \quad \text{(topics)}$$
$$z_{id} | \theta_d \sim \text{Mult}(\theta_d), \quad i \in |\mathbf{w}_d|, \quad \text{(topic indicators)}$$

## Come visit our poster 6pm Oren

| | | | |
| --- | --- | --- | --- |
| LDA | $87.8 \pm 0.6$ | $78.5 \pm 2.7$ | $66.3 \pm 2.6$ |
| movMF (tf) | $71.4 \pm 0.3$ | $64.5 \pm 0.6$ | $59.4 \pm 0.4$ |
| movMF (tf-idf) | $71.9 \pm 0.3$ | $74.2 \pm 0.4$ | $56.0 \pm 0.6$ |
| SAM (tf) | $88.6 \pm 0.4$ | $81.2 \pm 0.4$ | $70.5 \pm 0.5$ |
| SAM (tf-idf) | $93.3 \pm 0.3$ | $85.9 \pm 0.3$ | $75.0 \pm 0.4$ |
| **ds** | | | |
| | $91.8 \pm 0.4$ | $85.7 \pm 0.7$ | $75.6 \pm 0.8$ |
| | $91.1 \pm 0.3$ | $84.9 \pm 0.5$ | $75.8 \pm 0.8$ |
| | $91.4 \pm 0.5$ | $84.9 \pm 0.5$ | $75.3 \pm 0.6$ |
| | $91.9 \pm 0.4$ | $86.3 \pm 0.5$ | $75.6 \pm 0.6$ |
| | $94.1 \pm 0.3$ | $88.1 \pm 0.5$ | $78.1 \pm 0.6$ |

### Spherical Admixture Model

$$\boldsymbol{\mu}_t | \kappa_0 \sim \text{vMF}(\mathbf{m}, \kappa_0), \quad t \in T, \quad \text{(topic means)}$$
$$\phi_t | \boldsymbol{\mu}_t, \xi \sim \text{vMF}(\boldsymbol{\mu}_t, \xi), \quad t \in T, \quad \text{(topics)}$$
$$\boldsymbol{\beta}_d | \boldsymbol{\alpha} \sim \text{Dirichlet}(\boldsymbol{\alpha}), \quad d \in D, \quad \text{(topic proportions)}$$
$$\bar{\phi}_d | \phi, \boldsymbol{\beta}_d = \text{Avg}(\phi, \boldsymbol{\beta}_d), \quad d \in D, \quad \text{(spherical average)}$$
$$\mathbf{v}_d | \bar{\phi}_d, \kappa \sim \text{vMF}(\bar{\phi}_d, \kappa), \quad d \in D,$$

| | **Wikipedia** | | |
| --- | --- | --- | --- |
| $(+)$ | $(+)$ | $(-)$ | $(+)$ $(-)$ |
| navy | airport | album | opera | india | germany |
| ships | airlines | label | actor | temple | borough |
| naval | flights | singles | films | dynasty | england |
| submarines | bus | chart | players | indian | france |
| aircraft | satellites | song | conservatory | khan | parish |

1. Draw a set of $T$ topics $\phi$ on the unit h
2. For each document $d$, draw topic weig
3. Draw a document vector $\mathbf{v}_d$ from a vM

Top positive and negative term weights learned by SAM on the NIPS corpus and Wikipedia. (+) shows the highest weighted words and (−) shows lowest weighted within each topic.

Using SAM to generate features for document classification (L1 regularized logistic regression). Three different three-way classification tasks were derived from the 20-news dataset with increasing difficulty (i.e. classes become semantically similar).

Download: http://ww