# A Geometric Approach for Learning Latent Mixed Membership Models

Venkatesh Saligrama

Boston University, Boston, MA

(joint work with Weicong Ding & Prakash Ishwar)

# Outline

- ## Latent Mixture Models
    - Text Documents, User Preferences, Community Networks, …
    - Overall Goal/Objective: Algorithm with provable guarantees

- ## Topic Models & Estimation Problem

- ## Geometric Structure of Topic Models

- ## Algorithm & Guarantees: Exploiting Geometry

- ## Rank Aggregation Problem

- ## Real-World Expts.

Ding, Ishwar, Saligrama, ICML'13
Ding, Ishwar, Saligrama, NIPS'14
Ding, Ishwar, Saligrama, ITA'15
Ding, Ishwar, Saligrama, AISTATS'14
Ding, Ishwar, Saligrama, AISTATS'15

# Mixed membership latent variable model

- Text Docs ← (noisy) mixture of latent topics
- Connections in network ← mixture of latent communities
- User preferences ← mixtures of latent ranking factors

# Mixed membership latent variable model

Text document:

words    counts    Topics



**gene** — 3 — *Genetics*
**DNA** — 1
**genetic** — 2

**life** — 1
**evolve** — 1 — *Evolution*
**organism** — 2

**data** — 1
**number** — 1 — *Data Science*
**computer** — 3

…    …

**document = mixture of latent topics**

# Mixed membership latent variable model

words    counts    Influencing factors

User preferences



> 1   *"actor"*

> 0   *"music"*

> 1   *"special effect"*

**document = mixture of latent influencing factors**

# Overall Goal

- Learn/Estimate Latent Factors from Observations(docs)

- Goal: develop algorithms with

  - Provable guarantees
    - How many Docs to estimate Latent Factors within a tolerance?
    - Computational Cost: How does Algorithm scale with #params?

  ➔ *Model Fidelity*

  - Good empirical performance
    - Real-world datasets

  ➔ *Web Scale applications*

# Outline

- ## Latent Mixture Models
  - Text Documents, User Preferences, Community Networks, …

- ## Topic Models & Estimation Problem
  - Observation Model & Related Work

- ## Geometric Structure of Topic Models

- ## Algorithm & Guarantees: Exploiting Geometry

- ## Extensions to User Preference Model Estimation

- ## Empirical Results on Real-World Datasets

# "Bag of words" model: a text corpus example

One document in the collection:



| words | counts |
|---|---|
| **gene** | 3 |
| **DNA** | 1 |
| **genetic** | 2 |
| **life** | 1 |
| **evolve** | 1 |
| **organism** | 2 |
| **data** | 1 |
| **number** | 1 |
| **computer** | 3 |
| **…** | **…** |

$$
\begin{array}{c}
\text{word } 1 \\[6pt]
\text{word } 2 \\[30pt]
\\[30pt]
\\[30pt]
\text{word } W
\end{array}
\begin{bmatrix}
\beta_{11} & & \beta_{1K} \\
\beta_{21} & & \beta_{2K} \\
\vdots & \cdots & \vdots \\
\vdots & & \vdots \\
\beta_{W1} & & \beta_{WK}
\end{bmatrix}
$$

$$\text{topic } 1 \qquad \text{topic } K$$

Topic Matrix - $\beta$
- column = topic
- $W$ = vocabulary size
- $K$ = # topics

$$\begin{array}{c} \text{word 1} \\ \text{word 2} \\ \\ \\ \text{word } W \end{array} \begin{bmatrix} A_{11} & A_{12} & & A_{1M} \\ A_{21} & A_{22} & & A_{2M} \\ \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & & \vdots \\ A_{W1} & A_{W2} & & A_{WM} \end{bmatrix} = \begin{array}{c} \text{word 1} \\ \text{word 2} \\ \\ \\ \text{word } W \end{array} \begin{bmatrix} \beta_{11} & & \beta_{1K} \\ \beta_{21} & & \beta_{2K} \\ \vdots & \cdots & \vdots \\ \vdots & & \vdots \\ \beta_{W1} & & \beta_{WK} \end{bmatrix} \begin{bmatrix} \theta_{11} & \theta_{12} & & \theta_{1M} \\ \vdots & \vdots & \cdots & \vdots \\ \theta_{K1} & \theta_{K2} & & \theta_{KM} \end{bmatrix}$$

word 1, word 2 ... word $W$

doc. 1   doc. 2   doc. $M$

topic 1   topic $K$

doc. 1   doc. 2   doc. $M$

Document Distribution matrix - $A$
- column = distb. of a doc.
- $M$ = # docs.

Weight matrix - $\theta$
- column = mixing weights
- $M$ = # docs.

$$
\begin{array}{c}
\text{word 1} \\
\text{word 2} \\
\\
\\
\\
\text{word } W
\end{array}
\begin{bmatrix}
A_{11} & A_{12} & & A_{1M} \\
A_{21} & A_{22} & & A_{2M} \\
\vdots & \vdots & \cdots & \vdots \\
\vdots & \vdots & & \vdots \\
A_{W1} & A_{W2} & & A_{WM}
\end{bmatrix}
=
\begin{array}{c}
\text{word 1} \\
\text{word 2} \\
\\
\\
\text{word } W
\end{array}
\begin{bmatrix}
\beta_{11} & & \beta_{1K} \\
\beta_{21} & & \beta_{2K} \\
\vdots & \cdots & \vdots \\
\vdots & & \vdots \\
\beta_{W1} & & \beta_{WK}
\end{bmatrix}
\begin{bmatrix}
\theta_{11} & \theta_{12} & & \theta_{1M} \\
\vdots & \vdots & \cdots & \vdots \\
\theta_{K1} & \theta_{K2} & & \theta_{KM}
\end{bmatrix}
$$

$$
\begin{array}{ccc}
\text{doc.} & \text{doc.} & \text{doc.} \\
1 & 2 & M
\end{array}
\qquad
\begin{array}{cc}
\text{topic 1} & \text{topic } K
\end{array}
\qquad
\begin{array}{ccc}
\text{doc.} & \text{doc.} & \text{doc.} \\
1 & 2 & M
\end{array}
$$

$$\boxed{N \text{ iid samples + empirical words count}}$$

$$
\begin{array}{c}
\text{word 1} \\
\text{word 2} \\
\\
\\
\\
\text{word } W
\end{array}
\begin{bmatrix}
X_{11} & X_{12} & & X_{1M} \\
X_{21} & X_{22} & & X_{2Z} \\
\vdots & \vdots & \cdots & \vdots \\
\vdots & \vdots & & \vdots \\
X_{W1} & X_{W2} & & X_{WM}
\end{bmatrix}
$$

Observation matrix $X$
- column = word-freq. of a doc.
- $N$ = # word/doc.

$$
\begin{array}{c}
\text{word 1} \\
\text{word 2} \\
\\
\\
\\
\text{word } W
\end{array}
\begin{bmatrix}
A_{11} & A_{12} & & A_{1M} \\
A_{21} & A_{22} & & A_{2M} \\
\vdots & \vdots & \cdots & \vdots \\
\vdots & \vdots & & \vdots \\
A_{W1} & A_{W2} & & A_{WM}
\end{bmatrix}
=
\begin{array}{c}
\text{word 1} \\
\text{word 2} \\
\\
\\
\text{word } W
\end{array}
\begin{bmatrix}
\beta_{11} & & \beta_{1K} \\
\beta_{21} & & \beta_{2K} \\
\vdots & \cdots & \vdots \\
\vdots & & \vdots \\
\beta_{W1} & & \beta_{WK}
\end{bmatrix}
\begin{bmatrix}
\theta_{11} & \theta_{12} & & \theta_{1M} \\
\vdots & \vdots & \cdots & \vdots \\
\theta_{K1} & \theta_{K2} & & \theta_{KM}
\end{bmatrix}
$$

$$
\begin{array}{ccc}
\text{doc.} & \text{doc.} & \text{doc.} \\
1 & 2 & M
\end{array}
\qquad
\begin{array}{cc}
\text{topic 1} & \text{topic } K
\end{array}
\qquad
\begin{array}{ccc}
\text{doc.} & \text{doc.} & \text{doc.} \\
1 & 2 & M
\end{array}
$$

$$
\begin{array}{c}
\text{word 1} \\
\text{word 2} \\
\\
\\
\\
\text{word } W
\end{array}
\begin{bmatrix}
X_{11} & X_{12} & & X_{1M} \\
X_{21} & X_{22} & & X_{2Z} \\
\vdots & \vdots & \cdots & \vdots \\
\vdots & \vdots & & \vdots \\
X_{W1} & X_{W2} & & X_{WM}
\end{bmatrix}
$$

**Problem**
- Given : $X$ and $K$
- Goal : estimate $\beta$

| $K$ | # topics | ~100 |
|-----|----------|------|
| $W$ | vocab. size | ~10k |
| $N$ | #word/doc. | ~100 |
| $M$ | # doc. | ~100k |

# Related work

Topic matrix  Weight matrix

| Method | $\beta$ | $\theta$ | Approach | Issues |
|---|---|---|---|---|
| Nonnegative Matrix Factorization (NMF), e.g., [*Cichocki et al.,'09*] | Deterministic | Deterministic | Regularized joint optimization | NP Hard (Arora'12) Non-convex. Need approximations and heuristics. |
| "Bayesian Methods" e.g., LDA, CTM [*Blei et al.,'03*], | Deterministic or Prior | Prior | MAP or ML | Non-convex. Need approximations like MCMC. |
| Method of Moments [*Anandkumar et al.,'12,'13*] | Deterministic and sparse | Prior | Tensor decomposition | No empirical performance reported |
| Topic-separability based [*Several references*] | Approximate Separability | Prior | Geometric | |

# Outline

- ## Latent Mixture Models
    - Text Documents, User Preferences, Community Networks, …

- ## Topic Models & Estimation Problem
    - Related Work

- ## Geometric Structure of Topic Models
    - Inevitability of Separability in high-dimensions

- ## Algorithm & Guarantees: Exploiting Geometry

- ## Extensions to User Preference Model Estimation

- ## Empirical Results on Real-World Datasets

# Approximately Separable Topic Matrix ([Ding-Ishwar-S'14])



**Separable Topic Matrix**

λ-approximately separable if one word for each topic is predominantly unique

λ = 0 Case: Novel Word(s) **unique** to each topic

*[Boardman'93, Donoho'04]*, [Arora'13, Ding et. al.'13]

**Approximately Separable**

# Is Approximately Separable Fundamental?

| Dataset | W | K |
|---|---|---|
| Wikipedia | 109,611 | 50 |
| Twitter | 122,035 | 50 |
| New York Times | 102,660 | 100 |
| PubMed | 141,043 | 150 |

- **In real-world problems**
  - Size of vocab. $W$ >> #. Topics $K$

- **Main result:** Separability is an inevitable consequence of high-dimensionality!
  - Satisfied in estimates produced by NMF, LDA, and other algorithms (Bayesian Models)

**Generative Model**

# Why is Separability inevitable for W>> K?

- **Theorem:** Suppose $W \geq tKe^{\beta_0 K \log(K+1/\lambda)}$

$$\text{Prob}\{\beta \text{ not } \lambda - sep\} = O(W^{-t})$$

| Dataset | W | K |
|---|---|---|
| Wikipedia | 109,611 | 50 |
| Twitter | 122,035 | 50 |
| New York Times | 102,660 | 100 |
| PubMed | 141,043 | 150 |

Generative Model

# Separability in Practice

| Dataset | Vocab. size W | # Topics K | Prob. 0.01-separable |
|---|---|---|---|
| NIPS | 12,419 | 50 | 100±0% |
| Wikipedia | 109,611 | 50 | 99.9±0.3% |
| Twitter | 122,035 | 50 | 100±0% |
| New York Times | 102,660 | 100 | 99.6±0.6% |
| PubMed | 141,043 | 150 | 99.9±0.3% |

$\beta_0 = 0.01$, 1000 MC runs

**Generative Model**



- $\beta_0$ moderately small positive value in practice.
  - $\beta_0 = 0.01$ for $K \in [50, 200]$
- Some packages suggest $\beta_0 = c/W$ to get satisfactory empirical results.

$$W \geq tKe^{\beta_0 K \log(K+1/\lambda)}$$

- Analysis explains reasoning for this choice!

# Outline

- ## Latent Mixture Models
  - Text Documents, User Preferences, Community Networks, …

- ## Topic Models & Estimation Problem
  - Related Work

- ## Geometric Structure of Topic Models
  - Inevitability of Separability in high-dimensions

- ## Algorithm & Guarantees: Exploiting Geometry
  - Efficiently Identifying Extreme Points
  - Empirical Results on Real-World Datasets

- ## Extensions to User Preference Model Estimation

# Key Idea ($\lambda = 0$ case)

$$\boldsymbol{\beta}_{W \times K} \boldsymbol{\theta}_{K \times M} = \mathbf{A}_{W \times M}$$

|  | topic 1 | topic 2 | topic 3 |
|---|---|---|---|
| w1 | $\beta_1$ | 0 | 0 |
| w2 | $\beta_2$ | 0 | 0 |
| w3 | 0 | $\beta_3$ | 0 |
| w4 | 0 | $\beta_4$ | 0 |
| w5 | 0 | 0 | $\beta_5$ |
| w6 | 0 | 0 | $\beta_6$ |
|  | $\beta_{71}$ | $\beta_{72}$ | $\beta_{73}$ |
|  | $\cdots$ | | |

$$\begin{bmatrix} \leftarrow & \boldsymbol{\theta}_1 & \rightarrow \\ \leftarrow & \boldsymbol{\theta}_2 & \rightarrow \\ \leftarrow & \boldsymbol{\theta}_3 & \rightarrow \end{bmatrix}$$

**Weight Matrix**

$$\begin{bmatrix} \leftarrow & \beta_1\boldsymbol{\theta}_1 & \rightarrow \\ \leftarrow & \beta_2\boldsymbol{\theta}_1 & \rightarrow \\ \leftarrow & \beta_3\boldsymbol{\theta}_2 & \rightarrow \\ \leftarrow & \beta_4\boldsymbol{\theta}_2 & \rightarrow \\ \leftarrow & \beta_5\boldsymbol{\theta}_3 & \rightarrow \\ \leftarrow & \beta_6\boldsymbol{\theta}_3 & \rightarrow \\ \beta_{71}\boldsymbol{\theta}_1 + \beta_{72}\boldsymbol{\theta}_2 + \beta_{73}\boldsymbol{\theta}_3 \\ \cdots \end{bmatrix}$$

doc. 1 $\cdots$ doc. $M$

**Separable**
**Topic Matrix $\beta$**

**Distribution**
**Matrix $A$**



$\boldsymbol{\theta}_1$

$\boldsymbol{\theta}_2$

$\boldsymbol{\theta}_3$

$\Re_+^M$

# Key Idea ($\lambda = 0$ case)

$$\boldsymbol{\beta}_{W \times K} \boldsymbol{\theta}_{K \times M} = \mathbf{A}_{W \times M}$$

$$
\begin{array}{c}
\text{topic 1} \quad \text{topic 2} \quad \text{topic 3} \\
\text{w1} \begin{bmatrix} \beta_1 & 0 & 0 \\ \beta_2 & 0 & 0 \\ 0 & \beta_3 & 0 \\ 0 & \beta_4 & 0 \\ 0 & 0 & \beta_5 \\ 0 & 0 & \beta_6 \\ \beta_{71} & \beta_{72} & \beta_{73} \\ & \cdots & \end{bmatrix}
\end{array}
\qquad
\begin{array}{c}
\text{doc. 1} \quad \cdots \quad \text{doc. } M \\
\begin{bmatrix} \leftarrow & \boldsymbol{\theta}_1 & \rightarrow \\ \leftarrow & \boldsymbol{\theta}_2 & \rightarrow \\ \leftarrow & \boldsymbol{\theta}_3 & \rightarrow \end{bmatrix}
\end{array}
$$

**Weight Matrix**

$\longrightarrow$

$$
\begin{array}{c}
\text{doc. 1} \quad \cdots \quad \text{doc. } M \\
\begin{bmatrix} \leftarrow & \beta_1 \boldsymbol{\theta}_1 & \rightarrow \\ \leftarrow & \beta_2 \boldsymbol{\theta}_1 & \rightarrow \\ \leftarrow & \beta_3 \boldsymbol{\theta}_2 & \rightarrow \\ \leftarrow & \beta_4 \boldsymbol{\theta}_2 & \rightarrow \\ \leftarrow & \beta_5 \boldsymbol{\theta}_3 & \rightarrow \\ \leftarrow & \beta_6 \boldsymbol{\theta}_3 & \rightarrow \\ \beta_{71}\boldsymbol{\theta}_1 + \beta_{72}\boldsymbol{\theta}_2 + \beta_{73}\boldsymbol{\theta}_3 \\ \cdots \end{bmatrix}
\end{array}
$$

**Separable**
**Topic Matrix** $\beta$

$\boldsymbol{\theta}_1$, $\boldsymbol{\theta}_2$, $\boldsymbol{\theta}_3$, $\Re_+^M$

**Distribution**
**Matrix** $A$

# Key Idea

$$\widetilde{\beta}_{W \times K}\, \widetilde{\theta}_{K \times M} = \widetilde{A}_{W \times M}$$

doc. 1 ... doc. M

$$\begin{aligned}
&\leftarrow \quad \beta_1 \theta_1 \quad \rightarrow \\
&\leftarrow \quad \beta_2 \theta_1 \quad \rightarrow \\
&\leftarrow \\
&\leftarrow \quad \beta_4 \theta_2 \quad \rightarrow \\
&\leftarrow \quad \beta_5 \theta_3 \quad \rightarrow \\
&\leftarrow \quad \beta_6 \theta_3 \quad \rightarrow \\
&\beta_{71}\theta_1 + \beta_{72}\theta_2 + \beta_{73}\theta_3 \\
&\cdots
\end{aligned}$$

**Novel Word = Extreme Point**

**Novel word detection + Topic matrix estimation**

**Distribution Matrix $A$**

doc. 1 ... doc. M

$$\begin{aligned}
&\leftarrow \quad \widetilde{\theta}_1 \quad \rightarrow \\
&\leftarrow \quad \widetilde{\theta}_1 \quad \rightarrow \\
&\leftarrow \quad \widetilde{\theta}_2 \quad \rightarrow \\
&\leftarrow \quad \widetilde{\theta}_3 \quad \rightarrow \\
&\leftarrow \quad \widetilde{\theta}_3 \quad \rightarrow \\
&\widetilde{\beta}_{71}\widetilde{\theta}_1 + \widetilde{\beta}_{72}\widetilde{\theta}_2 + \widetilde{\beta}_{73}\widetilde{\theta}_3 \\
&\cdots
\end{aligned}$$

**Row Normalized Distribution Matrix $\widetilde{A}$**

$\widetilde{\theta}_1$

$\theta_1$

$\theta_2$

$\widetilde{\theta}_2$

$\theta_3$

$\widetilde{\theta}_3$

**Probability simplex** $\mathfrak{R}^M_+$

# Finite words/doc.

doc. 1 ... doc. M

$$\begin{bmatrix} \leftarrow & \beta_1\boldsymbol{\theta}_1 & \rightarrow \\ \leftarrow & \beta_2\boldsymbol{\theta}_1 & \rightarrow \\ \leftarrow & \beta_3\boldsymbol{\theta}_2 & \rightarrow \\ \leftarrow & \beta_4\boldsymbol{\theta}_2 & \rightarrow \\ \leftarrow & \beta_5\boldsymbol{\theta}_3 & \rightarrow \\ \leftarrow & \beta_6\boldsymbol{\theta}_3 & \rightarrow \\ \beta_{71}\boldsymbol{\theta}_1 + \beta_{72}\boldsymbol{\theta}_2 + \beta_{73}\boldsymbol{\theta}_3 \\ \cdots \end{bmatrix}$$

**Distribution Matrix $A$**

**row normalization** $\longrightarrow$

$\widetilde{\boldsymbol{\theta}}_1$

$\widetilde{\boldsymbol{\theta}}_2$

$\widetilde{\boldsymbol{\theta}}_3$

**Probability simplex**

doc. 1 ... doc. M

$$\begin{bmatrix} \leftarrow & \widetilde{\boldsymbol{\theta}}_1 & \rightarrow \\ \leftarrow & \widetilde{\boldsymbol{\theta}}_1 & \rightarrow \\ \leftarrow & \widetilde{\boldsymbol{\theta}}_2 & \rightarrow \\ \leftarrow & \widetilde{\boldsymbol{\theta}}_2 & \rightarrow \\ \leftarrow & \widetilde{\boldsymbol{\theta}}_3 & \rightarrow \\ \leftarrow & \widetilde{\boldsymbol{\theta}}_3 & \rightarrow \\ \widetilde{\beta}_{71}\widetilde{\boldsymbol{\theta}}_1 + \widetilde{\beta}_{72}\widetilde{\boldsymbol{\theta}}_2 + \widetilde{\beta}_{73}\widetilde{\boldsymbol{\theta}}_3 \\ \cdots \end{bmatrix}$$

**Row Normalized Distribution Matrix $A$**

# Finite words/doc. $\widetilde{\boldsymbol{\beta}}_{W \times K} \widetilde{\boldsymbol{\theta}}_{K \times M} = \widetilde{\mathbf{A}}_{W \times M} \approx \widetilde{\mathbf{X}}_{W \times M}$

**Key issue:**
$N$ fixed ➡ perturbation
does not vanish

$$\begin{array}{ccc}
\text{doc.} & \cdots & \text{doc.} \\
1 & & M
\end{array}$$

$$\begin{bmatrix}
\leftarrow & \widetilde{\boldsymbol{\theta}}_1 & \rightarrow \\
\leftarrow & \widetilde{\boldsymbol{\theta}}_1 & \rightarrow \\
\leftarrow & \widetilde{\boldsymbol{\theta}}_2 & \rightarrow \\
\leftarrow & \widetilde{\boldsymbol{\theta}}_2 & \rightarrow \\
\leftarrow & \widetilde{\boldsymbol{\theta}}_3 & \rightarrow \\
\leftarrow & \widetilde{\boldsymbol{\theta}}_3 & \rightarrow \\
\widetilde{\beta}_{71}\widetilde{\boldsymbol{\theta}}_1 + \widetilde{\beta}_{72}\widetilde{\boldsymbol{\theta}}_2 + \widetilde{\beta}_{73}\widetilde{\boldsymbol{\theta}}_3 \\
\cdots
\end{bmatrix}$$

**Row Normalized Distribution Matrix $\widetilde{A}$**

| | | |
|---|---|---|
| $K$ | # topics | ~100 |
| $W$ | vocab. size | ~10k |
| $N$ | #word/doc. | ~100 |
| $M$ | # doc. | ~100k |

**Probability simplex**

$$\begin{array}{ccc}
\text{doc.} & \cdots & \text{doc.} \\
1 & & M
\end{array}$$

$$\begin{bmatrix}
\leftarrow & \widetilde{X}_1 & \rightarrow \\
\leftarrow & \widetilde{X}_2 & \rightarrow \\
\leftarrow & \widetilde{X}_3 & \rightarrow \\
\leftarrow & \widetilde{X}_4 & \rightarrow \\
\leftarrow & \widetilde{X}_5 & \rightarrow \\
\leftarrow & \widetilde{X}_6 & \rightarrow \\
\leftarrow & \widetilde{X}_7 & \rightarrow \\
\cdots
\end{bmatrix}$$

**Row Normalized Observation Matrix: $\widetilde{X}$**

**Nonnegative**

$$\widetilde{\beta}_{W \times K} \overbrace{\widetilde{\theta}_{K \times M}} = \widetilde{A}_{W \times M}$$

**Separable**
**Nonnegative**

| $K$ | # topics | ~100 |
|-----|----------|------|
| $W$ | vocab. size | ~10k |
| $N$ | #word/doc. | ~100 |
| $M$ | # doc. | ~100k |

**Nonnegative**

$$M\widetilde{X}\widetilde{X}^T \xrightarrow[M \to \infty]{a.s} \widetilde{\beta}(M\widetilde{\theta}\ \widetilde{\theta}^T)\widetilde{\beta}^T = E_{W \times W}$$

**Separable**
**Nonnegative**

$\widetilde{\theta}_1 = \widetilde{A}_1 = \widetilde{A}_2$

**Probability simplex**

$\widetilde{\theta}_2 = \widetilde{A}_3 = \widetilde{A}_4$

$\widetilde{\theta}_3 = \widetilde{A}_5 = \widetilde{A}_6$

$E_1 = E_2$

$\Re_+^W$

$E_3 = E_4$

$E_5 = E_6$

**Novel word detection + Topic matrix estimation**

# Detect Novel Words via Projections

- Max/Min of projection → extreme points of convex hull

- Which directions to use
  →Generate $P$ iid Isotropy directions

# "Robustness" of extreme points
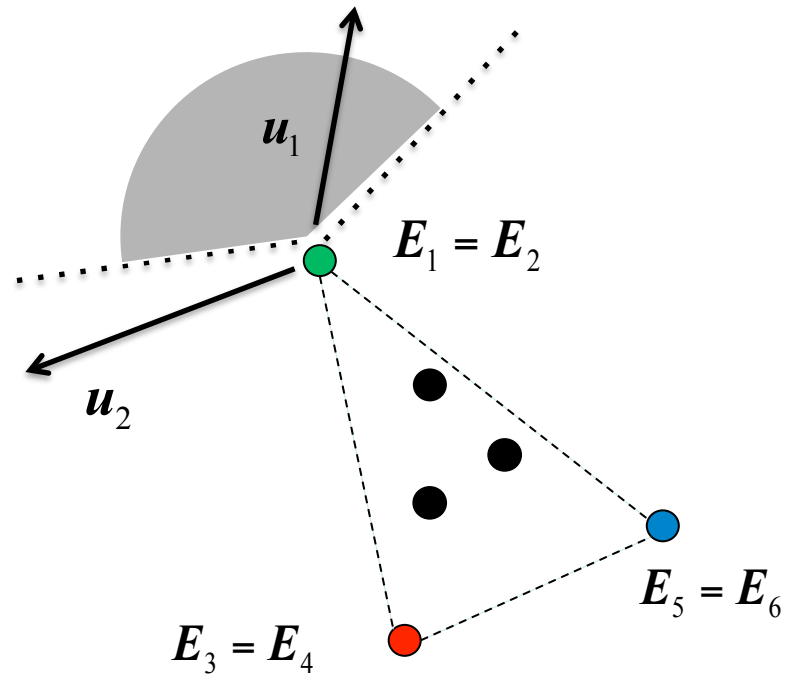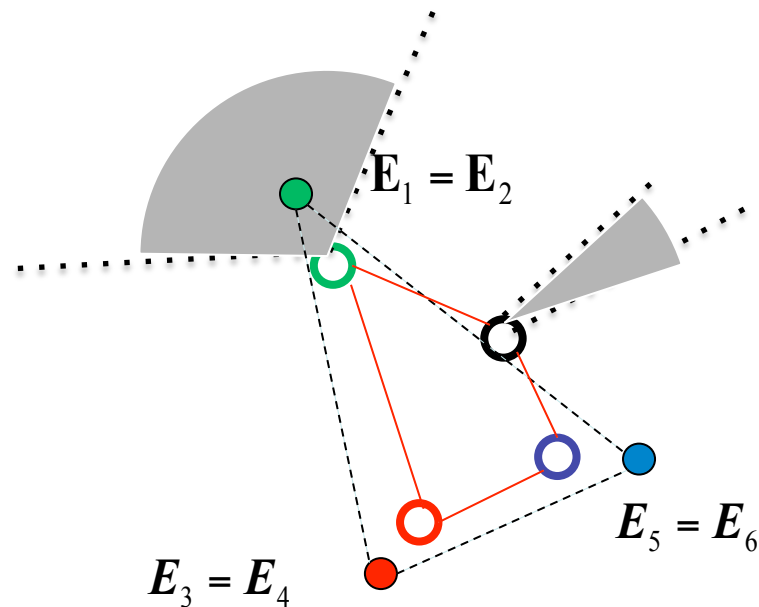
- Max/Min of projection
  → extreme points of convex hull

- Which directions to use
  →Generate a few iid Isotropy directions

- Freq. of maximum ≈ Solid Angle of an Extreme Point

# Extension to Approximately Separable ($\lambda > 0$)

- Approx. novel words ⇔ larger solid angles

- Solution
  - → Sort solid angles
  - → Take the top – K extreme points

# Main result [Ding et al, '13,'14]

| $K$ | # topics | ~100 |
|-----|----------|------|
| $W$ | vocab. size | ~10k |
| $N$ | #word/doc. | ~100 |
| $M$ | # doc. | ~100k |

- Computational complexity :

$$O(MNK + WK + WK^3)$$

- Sample complexity :

Under the **Simplicial Condition** on $R$', with $u \sim N(0, I_W)$, the proposed Random Projection algorithm can detect all novel words of $K$ topics with probability $1\text{-}\delta$ if

$$M \geq \text{Poly}\left(W, \log\left(\frac{1}{\delta}\right), K, \frac{1}{N}\right), P \geq \text{Poly}\left(\log W, \log\left(\frac{1}{\delta}\right), K\right)$$

Moreover, if R is **full-rank**, can recover $\beta$ with $\varepsilon$ element-wise error with probability $1\text{-}\delta$.

# Distributed Implementation



**Server 2**

**Server 1**

Processor 1

Memory 1

Processor 2

Memory 2

$\mathbf{X}_{(2)}$

**Server** $L$

Processor $L$

Memory $L$

$\mathbf{X}_{(1)}$

Processor

Memory

**Fusion Center**

$\mathbf{X}_{(L)}$

$$\hat{\mathbf{E}}\mathbf{d}_r = \sum_{l=1}^{L} \mathbf{X}_{(l)} \mathbf{X}_{(l)}^{T} \mathbf{d}_r$$

- Modern web-scale database are distributed
- $M$ document archived on $L$ servers
- Goal: low communication cost O(WK)

# Experimental Results (semi-synthetic data)

Real-world corpus
New York Times articles

→

Topic matrix learnt
by Gibbs sampling

Generate synthetic
docs. with Dirichlet
prior

Add artificial novel
words;
Generate synthetic
docs.

| $M$ | 300,000 |
|---|---|
| $N$ | 300 |
| $W$ | 14943 |
| $K$ | 100 |
| $L$ | 200 |

- Semi-synthetic data can resemble the dimensionality and sparsity of real-world data

# Experimental Results (semi-synthetic data)

# Experimental Results (semi-synthetic data)

| | Semi-syn +novel NYT | Semi-syn NYT |
|---|---|---|
| $M$ | Variable | Variable |
| $N$ | 300 | 300 |
| $W$ | 15043 words, 100 novel | 14943 |
| $K$ | 100 | 100 |
| $L$ | 200 | 200 |

**\* T (Gibbs) ~ 6918 sec**



**RP**

**RecoverL2**

**DDP**
[*Ding et al.,13'*]

*Semi-synthetic +Novel NYT*

*l1 error*

**Comp. Time, in sec.**

# *Experimental* Results (Real-World Text Corpus)

New York Times dataset

| $M$ | 300,000 = 240k train + 60k test |
|---|---|
| $N$ | 300 words/doc. (avg.) |
| $W$ | 14,943 |
| $K$ | 50/100/150 topics |

Decreasing Prob.

| **"Weather"** | **"Emotion"** | **"Politics"** | **"Football"** |
|---|---|---|---|
| Weather | Feeling | Election | Yard |
| Wind | Sense | *Florida* | Game |
| Air | Love | Ballot | Team |
| Storm | Character | Vote | Season |
| Rain | Heart | *Al_gore* | Play |
| Cold | Emotion | Recount | *NFL* |

(See Ding et al., '13 for more example topics)

# Outline

- ## Latent Mixture Models
  - Text Documents, User Preferences, Community Networks, …

- ## Topic Models & Estimation Problem
  - Related Work

- ## Geometric Structure of Topic Models
  - Inevitability of Separability in high-dimensions

- ## Algorithm & Guarantees: Exploiting Geometry
  - Efficiently Identifying Extreme Points
  - Empirical Results on Real-World Datasets

- ## Rank Aggregation Problem
  - Heterogenous Population ~ Mixture of Mallows Model
  - Reduction to Topic Modeling Problem
  - Empirical Results on Real-World Datasets

# Mixed membership latent variable model

User preferences

words counts

Influencing factors



> 1 *"actor"*

> 0 *"music"*

> 1 *"special effect"*

**document = mixture of latent influence factors**

actor          musical

Prob. Prefer movie 1 over 2 in 1<sup>st</sup> latent factor

$(1,2)$

$(1,3)$

$(Q\text{-}1.Q)$

$$\begin{bmatrix} \beta_{11} & & \beta_{1K} \\ \beta_{21} & & \beta_{2K} \\ \vdots & \cdots & \vdots \\ \vdots & & \vdots \\ \beta_{W1} & & \beta_{WK} \end{bmatrix}$$

topic 1        topic $K$

Ranking matrix - $\boldsymbol{\beta}$
- column = "topic"
- $W$ = # ordered pairs
      = $Q\,(Q - 1)$
- $K$ = # topics

- Generative Model for Latent factor

- Mallows model
  - Baseline permutation $\sigma_0$
  - Prob. of permutation $\sigma$:    $Prob\{\sigma \mid \sigma_0\} \propto \phi^{dist(\sigma, \sigma_0)}$ $\implies \beta$
  - Heterogeneous population
    - Dispersion factor $\phi$

# Mixed Membership Rank Aggregation Problem

"Plate" representation

Key parameters



| $Q$ | # items |
|---|---|
| $W=Q(Q\text{-}1)$ | # ordered pairs |
| $K$ | # latent rankings |
| $M$ | # users |
| $N$ | # comps./user |

➔ Words

➔ Topics

➔ Documents

$w$: comparisons $(i,j)$

# Related Work

| Category | Models | Issues | Hetero-geneity | User-inconsistency |
|---|---|---|---|---|
| Single ranking models | Mallows[*Mallows,'57*], BTL [*Bradley & Terry,'52*], etc. | Only one global ranking | Only via deviation from base ranking scheme | Via noisy observation model |
| Mixture of ranking models | Mixture of Mallows model [*Lu & Boutilier ICML11, Awasthi et al. NIPS14*], Mixture of single rankings [*Farias et al. NIPS09*], Mixture of BTLs [*Oh & Shah NIPS14*] | Each user is dominated by one type | Multiple mixture components | Via noisy observation model |
| "Topic" ranking models | [*Ding et al. NIPS14*], Mixed membership Mallows, etc. | | Multiple shared rankings | Via probabilistic mixture + noisy observation |

# Approximate Separability

|  | Rank 1 | Rank 2 | Rank 3 |
|---|---|---|---|
| Pair 1 | 0.98 | 0.01 | 0.01 |
| Pair 2 | 0.01 | 0.99 | 0.01 |
| Pair 3 | 0.01 | 0.01 | 0.90 |
| Pair 4 | 0.98 | 0.90 | 0.10 |
| Pair 5 | 0.10 | 0.09 | 0.90 |
|  |  | ... |  |

$\boldsymbol{\beta}$

- Most ranking matrix are $\lambda$-Approximate separable, # items $Q >>$ # factors $K$

$$\Pr(\boldsymbol{\sigma} \text{ is } \lambda\text{ - separable}) \geq 1 - K \exp(-QL(\lambda;\phi)^{-2K+1})$$

Approximately Separable ranking matrix

| $\phi$ | Prob. of 0.05-separable |
|---|---|
| 0.0 | 93.3% |
| 0.1 | 87.0% |
| 0.2 | 79.3% |
| 0.5 | 42.6% |

$Q = 100$

$K = 10$

1000 Monte Carlo runs

| $Q$ | # items |
|---|---|
| $W$ | # ordered pairs |
| $K$ | # latent ranking |
| $M$ | # users |
| $N$ | # comp. / user |
| $P$ | # projections |

# Main result

- Computational complexity :

$$O(MNP + Q^2 K^3)$$

- Sample complexity :

- If
  - Correlation matrix of weight prior has **full-rank** and
  - The ranking matrix **σ** is $\lambda$-**separable** and
  - $\lambda \leq cK^{-2}$
- Then
  - Proposed Random Projection algorithm can estimate the ranking matrix correctly with probability at least $1\text{-}\delta$ for all
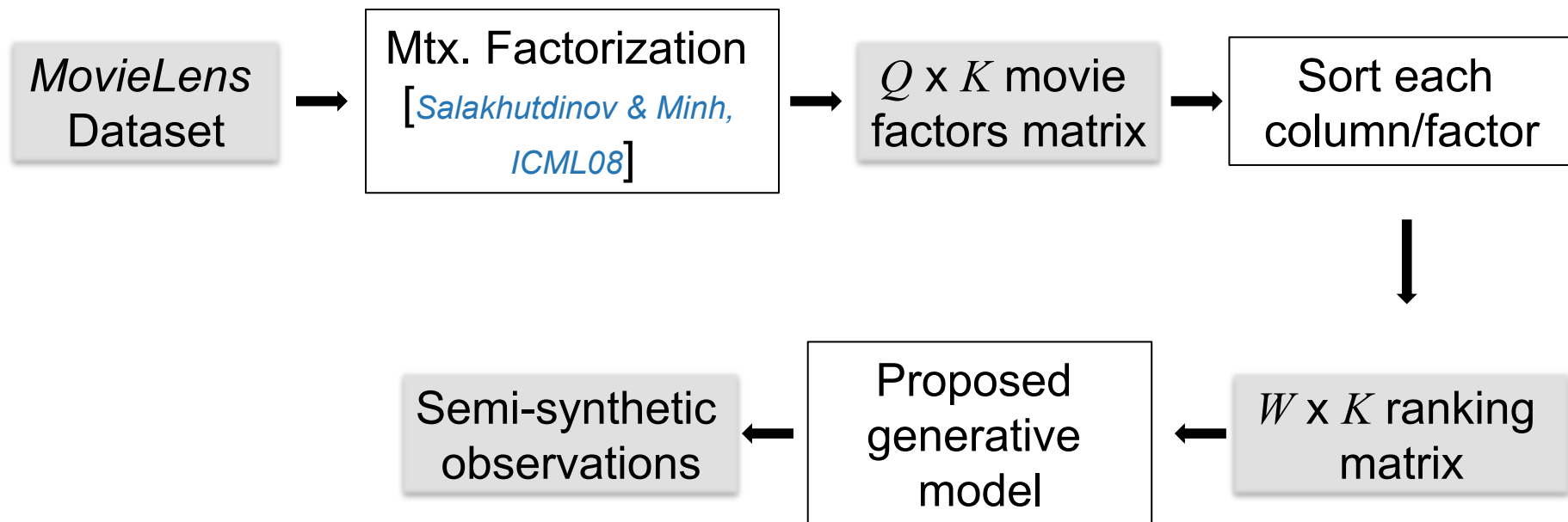
$$M \geq \mathrm{Poly}\left(W, K, \frac{1}{N}, \log\left(\frac{1}{\delta}\right)\right), \ P \geq \mathrm{Poly}\left(\log W, \log\left(\frac{1}{\delta}\right), K\right)$$

# Semi-synthetic data

| $Q$ | 100 most rated movies |
|-----|------------------------|
| $K$ | 10 latent rankings |
| $M$ | 5940 users |
| | ~200K ratings |

- To resemble the dimensionality and characteristics of real-world data

*MovieLens* Dataset → Mtx. Factorization [*Salakhutdinov & Minh, ICML08*] → $Q$ x $K$ movie factors matrix → Sort each column/factor

→ $W$ x $K$ ranking matrix → Proposed generative model → Semi-synthetic observations

# Semi-synthetic data



| $Q$ | 100 most rated movies |
|-----|------------------------|
| $K$ | 10 latent rankings |
| $N$ | 300 comparisons/user |
| $M$ | # user, variable |

- Dirichlet Prior for $\boldsymbol{\theta}$
- Uniform distribution for $\boldsymbol{\mu}$
- $\phi = \phi_1 = \ldots = \phi_K$

# *Movielens* dataset – new comparison

- **Data:** generate comparisons from ratings

  User 1: Movie A, 4 star Movie C, 2 star ➜ $w = (A, C)$ for user 1

- **Task**: predict new comparisons for users in the training set

- **Measure**: predictive log-probability [*Wallach et al. ICML09*]

$$\frac{1}{\text{Total \# test pairs}} \sum_{i,m} \log\{p(i \text{ th test pair} \mid \hat{\sigma}, \text{Training pairs of user } m)\}$$

| | |
|---|---|
| $Q$ | 100 most rated movies |
| $K$ | variable |
| $M$ | 5940 users |



RP

FJS
[*Farias et al. NIPS09*]

# *Movielens* dataset – new user

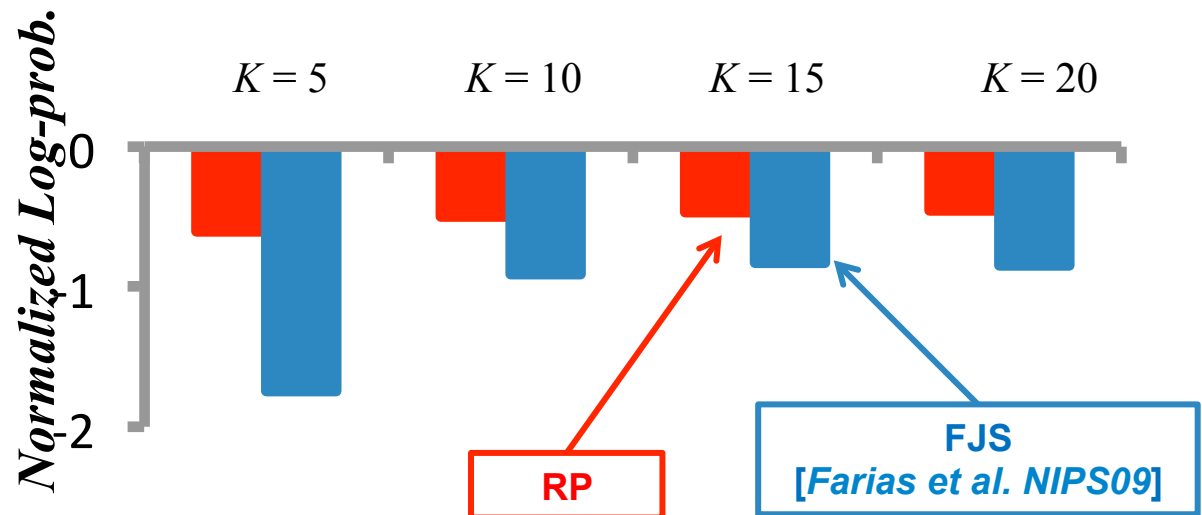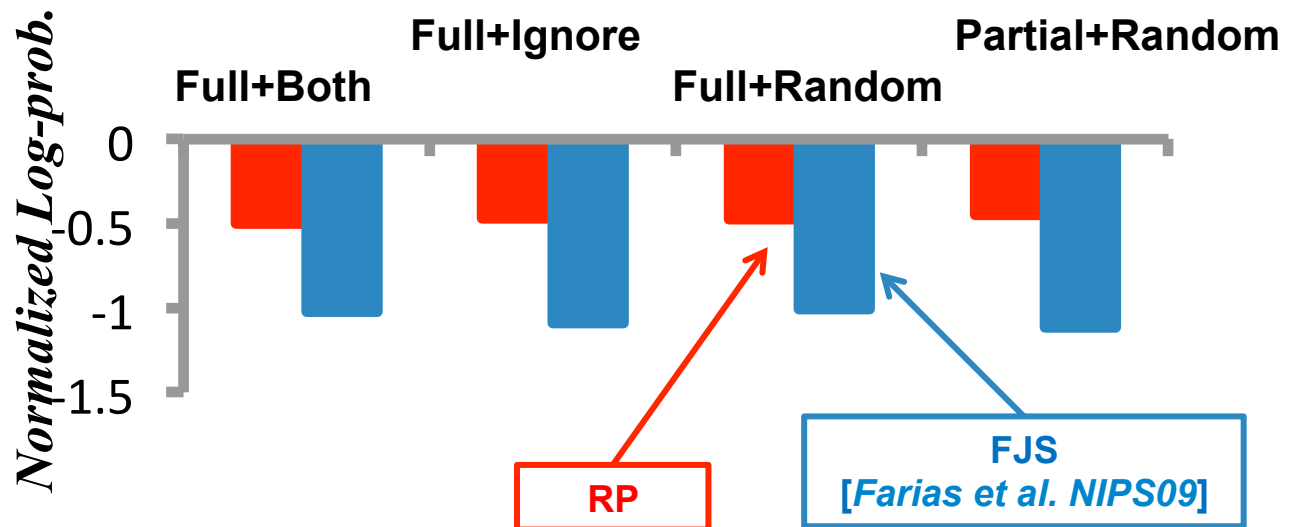- **Data:** generate comparisons from ratings

    User 1: Movie A, 4 star Movie C, 2 star➔ $w = (A, C)$ for user 1

- **Task**: predict comparison of new users

- **Measure**: predictive log-probability [*Wallach et al. ICML09*]

$$\frac{1}{\text{Total \# test pairs}} \sum\nolimits_{i,m} \log\{p(i \text{ th test pair of user } m \mid \hat{\boldsymbol\sigma})\}$$

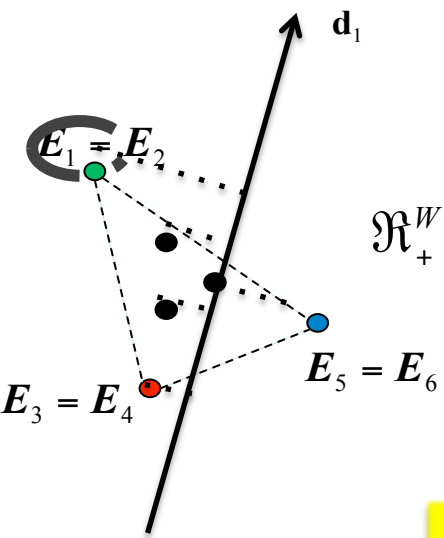| $Q$ | 100 most rated movies |
|---|---|
| $K$ | 10 latent rankings |
| $M$ | 4000 *training* user |
| $M$ | 2040 *testing* user |

# *Movielens* dataset – predicting stars

- Predict star ratings via ranking models
    - Generate comparisons from training ratings
    - Learn mixed membership Mallows model with Dirichlet prior
    - For each testing movie review:

- Measure: RMSE of estimated star ratings

| K | PMF | BPMF | BPMF-int | TM(Ding et al.,14) | MMMM |
|---|------|-------|----------|-----------------|-------|
| 10 | 1.0491 | 0.8254 | 0.8723 | 0.8840 | 0.8509 |
| 15 | 0.9127 | 0.8236 | 0.8734 | 0.8780 | 0.8296 |
| 20 | 0.9250 | 0.8213 | 0.8678 | 0.8721 | 0.8241 |

**Rating based models**

# Summary

|  | Rank 1 | Rank 2 | Rank 3 |
|---|---|---|---|
| Pair 1 | 0.98 | 0.01 | 0.01 |
| Pair 2 | 0.01 | 0.99 | 0.01 |
| Pair 3 | 0.01 | 0.01 | 0.90 |
| Pair 4 | 0.98 | 0.90 | 0.10 |
| Pair 5 | 0.10 | 0.09 | 0.90 |
|  | … | | |

$$\beta$$

High-D Latent Factor Models
Geometry ~ Approx Sep.

Simple geometric picture

Efficient randomized algorithm

Consistency, efficiency, state-of-the-art performance

Approximately Separable
ranking matrix

$d_1$

$E_1 = E_2$

$\Re_+^W$

$E_3 = E_4$

$E_5 = E_6$