

Information Retrieval and Topic Models

Mausam

(Based on slides of W. Arms, Dan Jurafsky, Thomas Hofmann,
Ata Kaban, Chris Manning, Melanie Martin)

Unstructured data in 1620

- Which plays of Shakespeare contain the words ***Brutus*** ***AND Caesar*** but ***NOT Calpurnia***?
- One could grep all of Shakespeare's plays for ***Brutus*** and ***Caesar***, then strip out lines containing ***Calpurnia***?
- Why is that not the answer?
 - Slow (for large corpora)
 - ***NOT Calpurnia*** is non-trivial
 - Other operations (e.g., find the word ***Romans*** near ***countrymen***) not feasible
 - Ranked retrieval (best documents to return)
 - Later lectures

Term-document incidence matrices

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0

Brutus AND Caesar BUT NOT Calpurnia

1 if play contains
word, 0 otherwise

Incidence vectors

- So we have a 0/1 vector for each term.
- To answer query: take the vectors for ***Brutus***, ***Caesar*** and ***Calpurnia*** (complemented) → bitwise ***AND***.
 - 110100 ***AND***
 - 110111 ***AND***
 - 101111 =
 - **100100**

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0

Answers to query

• Antony and Cleopatra, Act III, Scene ii

Agrippa [Aside to DOMITIUS ENOBARBUS]: Why, Enobarbus,
When Antony found Julius **Caesar** dead,
He cried almost to roaring; and he wept
When at Philippi he found **Brutus** slain.

• Hamlet, Act III, Scene ii

Lord Polonius: I did enact Julius **Caesar** I was killed i' the
Capitol; **Brutus** killed me.



Term-document count matrices

- Consider the number of occurrences of a term in a document:
 - Each document is a count vector in $\mathbb{N}^{|V|}$: a column below

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0

tf-idf weighting

- The tf-idf weight of a term is the product of its tf weight and its idf weight.

$$w_{t,d} = (1 + \log \text{tf}_{t,d}) \times \log_{10}(N / \text{df}_t)$$

- **Best known weighting scheme in information retrieval**
 - Note: the “-” in tf-idf is a hyphen, not a minus sign!
 - **Alternative names: tf.idf, tf x idf**
- Increases with the number of occurrences within a document
- **Increases with the rarity of the term in the collection**

Final ranking of documents for a query

$$\text{Score}(q, d) = \sum_{t \in q \cap d} \text{tf.idf}_{t,d}$$

Binary \rightarrow count \rightarrow weight matrix

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	5.25	3.18	0	0	0	0.35
Brutus	1.21	6.1	0	1	0	0
Caesar	8.59	2.54	0	1.51	0.25	0
Calpurnia	0	1.54	0	0	0	0
Cleopatra	2.85	0	0	0	0	0
mercy	1.51	0	1.9	0.12	5.25	0.88
worser	1.37	0	0.11	4.15	0.25	1.95

Each document is now represented by a real-valued vector of tf-idf weights $\in \mathbb{R}^{|V|}$

Documents as vectors

- Now we have a $|V|$ -dimensional vector space
- Terms are axes of the space
- Documents are points or vectors in this space
- Very high-dimensional: tens of millions of dimensions when you apply this to a web search engine
- These are very sparse vectors – most entries are zero

Queries as vectors

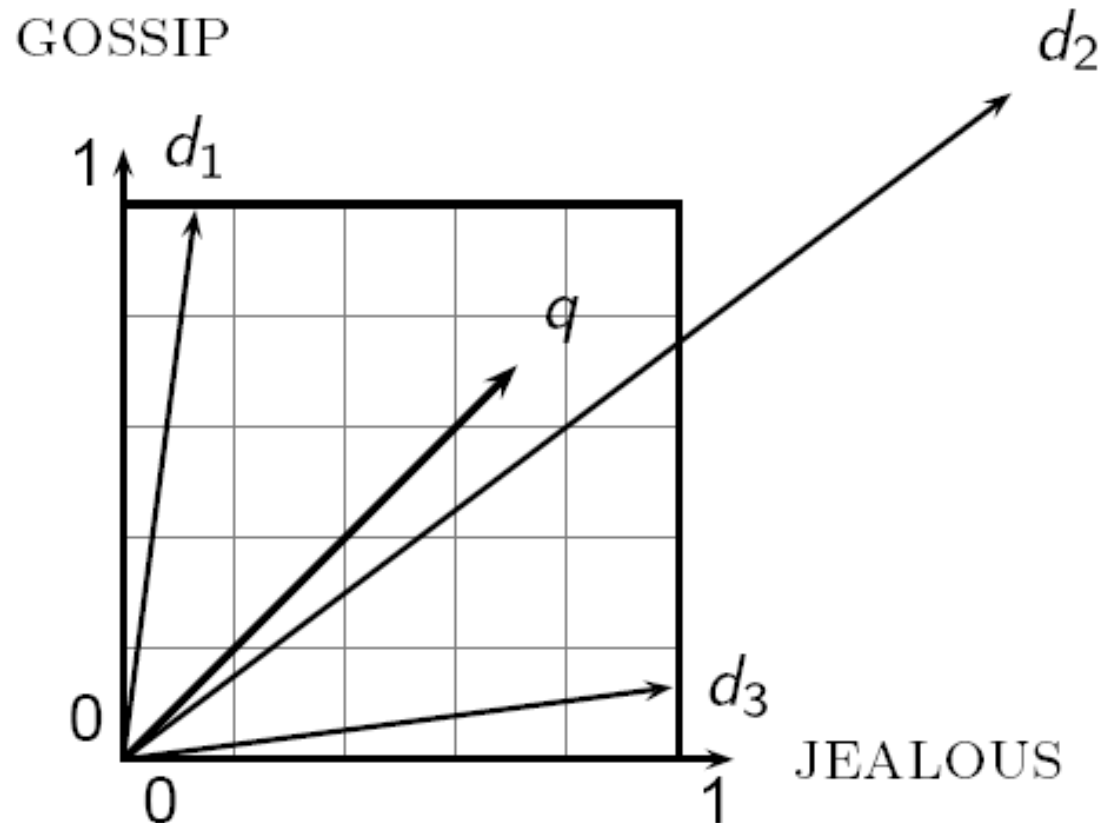
- Key idea 1: Do the same for queries: represent them as vectors in the space
- Key idea 2: Rank documents according to their proximity to the query in this space
- proximity = similarity of vectors
- proximity \approx inverse of distance
- **Recall: We do this because we want to get away from the you're-either-in-or-out Boolean model**
- Instead: rank more relevant documents higher than less relevant documents

Formalizing vector space proximity

- First cut: distance between two points
 - (= distance between the end points of the two vectors)
- **Euclidean distance?**
- Euclidean distance is a bad idea . . .
- . . . because Euclidean distance is **large** for vectors of **different lengths**.

Why distance is a bad idea

The Euclidean distance between q and \vec{d}_2 is large even though the distribution of terms in the query \vec{q} and the distribution of terms in the document \vec{d}_2 are very similar.



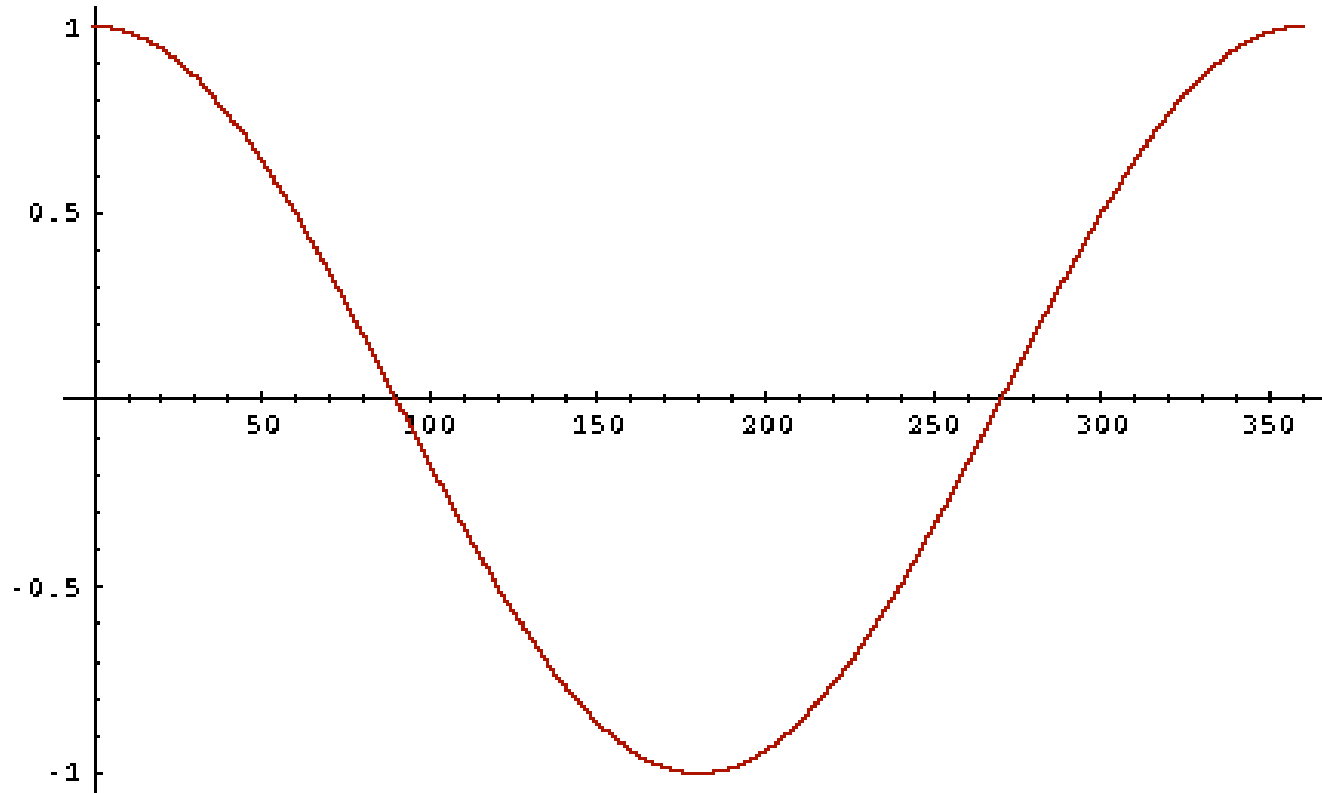
Use angle instead of distance

- Thought experiment: take a document d and append it to itself. Call this document d' .
- “Semantically” d and d' have the same content
- The Euclidean distance between the two documents can be quite large
- The angle between the two documents is 0, corresponding to maximal similarity.
- Key idea: Rank documents according to angle with query.

From angles to cosines

- The following two notions are equivalent.
 - Rank documents in decreasing order of the angle between query and document
 - Rank documents in increasing order of $\text{cosine}(\text{query}, \text{document})$
- Cosine is a monotonically decreasing function for the interval $[0^\circ, 180^\circ]$

From angles to cosines



- But how – *and why* – should we be computing cosines?

Length normalization

- A vector can be (length-) normalized by dividing each of its components by its length – for this we use the L_2 norm:
$$\|x\|_2 = \sqrt{\sum_i x_i^2}$$
- Dividing a vector by its L_2 norm makes it a unit (length) vector (on surface of unit hypersphere)
- Effect on the two documents d and d' (d appended to itself) from earlier slide: they have identical vectors after length-normalization.
 - Long and short documents now have comparable weights

cosine(query,document)

$$\cos(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}| |\vec{d}|} = \frac{\vec{q}}{|\vec{q}|} \cdot \frac{\vec{d}}{|\vec{d}|} = \frac{\sum_{i=1}^{|\mathcal{V}|} q_i d_i}{\sqrt{\sum_{i=1}^{|\mathcal{V}|} q_i^2} \sqrt{\sum_{i=1}^{|\mathcal{V}|} d_i^2}}$$

Dot product
Unit vectors

q_i is the tf-idf weight of term i in the query

d_i is the tf-idf weight of term i in the document

$\cos(\vec{q}, \vec{d})$ is the cosine similarity of \vec{q} and \vec{d} ... or, equivalently, the cosine of the angle between \vec{q} and \vec{d} .

Similarity Measures Compared

$$|Q \cap D|$$

Simple matching (coordination level match)

$$2 \frac{|Q \cap D|}{|Q| + |D|}$$

Dice's Coefficient

$$\frac{|Q \cap D|}{|Q \cup D|}$$

Jaccard's Coefficient

$$\frac{|Q \cap D|}{|Q|^{1/2} \times |D|^{1/2}}$$

Cosine Coefficient (what we studied)

$$\frac{|Q \cap D|}{\min(|Q|, |D|)}$$

Overlap Coefficient

Summary – vector space ranking

- Represent the query as a weighted tf-idf vector
- Represent each document as a weighted tf-idf vector
- Compute the cosine similarity score for the query vector and each document vector
- Rank documents with respect to the query by score
- Return the top K (e.g., $K = 10$) to the user

Evaluating ranked results: Mean Avg Precision

- 1 R
- 2 N
- 3 N
- 4 R
- 5 R
- 6 N
- 7 R
- 8 N
- 9 N
- 10 N

Assume 10 rel docs
in collection

Common evaluation measure...

- Mean average precision (MAP)
 - AP: Average of the precision value obtained for the top k documents, each time a relevant doc is retrieved
 - Avoids interpolation, use of fixed recall levels
 - Does weight most accuracy of top returned results
 - MAP for set of queries is arithmetic average of APs
 - Macro-averaging: each query counts equally

Problems



- Synonyms: separate words that have the same meaning.
 - E.g. 'car' & 'automobile'
 - They tend to reduce recall
 - Polysems: words with multiple meanings
 - E.g. 'Java'
 - They tend to reduce precision
- The problem is more general: there is a disconnect between topics and words

- ‘... a more appropriate model should consider some *conceptual* dimensions instead of words.’ (Gardenfors)

Latent Semantic Analysis (LSA)

- LSA aims to discover something about the meaning behind the words; about the topics in the documents.
- What is the difference between topics and words?
 - Words are observable
 - Topics are not. They are latent.
- How to find out topics from the words in an automatic way?
 - We can imagine them as a compression of words
 - A combination of words
 - Try to formalise this

Latent Semantic Analysis

- Singular Value Decomposition (SVD)
 - $A(m*n) = U(m*r) E(r*r) V(r*n)$
 - Keep only k eigen values from E
 - $A(m*n) = U(m*k) E(k*k) V(k*n)$
 - Convert terms and documents to points in k -dimensional space
 - Low-rank approximation

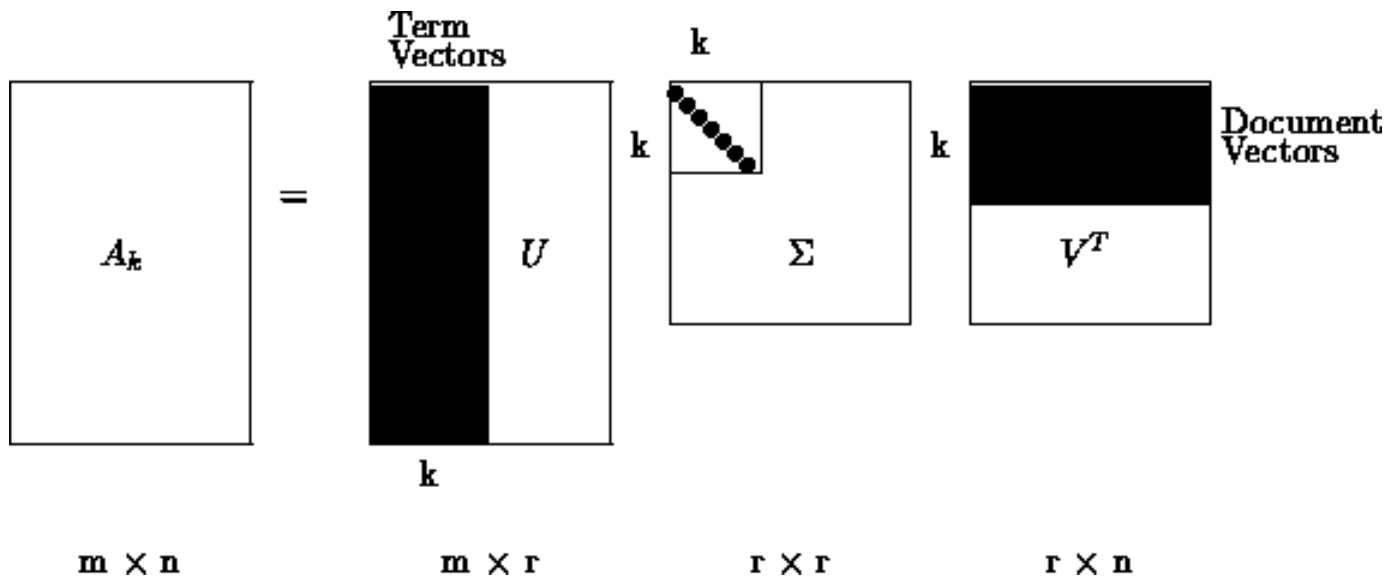
Latent Semantic Analysis

- Singular Value Decomposition

$$\{A\} = \{U\}\{S\}\{V\}^T$$

- Dimension Reduction

$$\{\tilde{A}\} \sim \{\tilde{U}\}\{\tilde{S}\}\{\tilde{V}\}^T$$



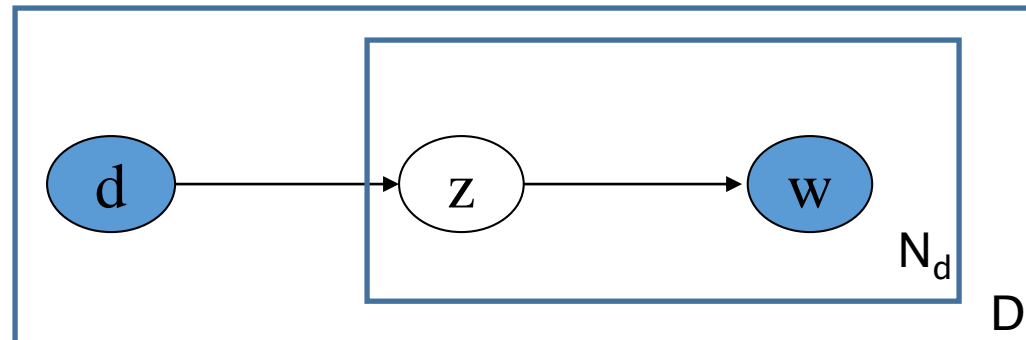
Latent Semantic Analysis

- LSA puts documents together even if they don't have common words if
 - The docs share frequently co-occurring terms
- Disadvantages:
 - Statistical foundation is missing

PLSA addresses this concern!

PLSA

- Latent Variable model for general co-occurrence data
 - Associate each observation (w,d) with a class variable $z \in Z\{z_1, \dots, z_K\}$
- Generative Model
 - Select a doc with probability $P(d)$
 - Pick a latent class z with probability $P(z | d)$
 - Generate a word w with probability $p(w | z)$



PLSA

- To get the joint probability model

$$P(d, w) = P(d)P(w|d), \text{ where}$$

$$P(w|d) = \sum_{z \in \mathcal{Z}} P(w|z)P(z|d) .$$

$$P(d, w) = \sum_{z \in \mathcal{Z}} P(z)P(w|z)P(d|z).$$

Model fitting with EM

- We have the equation for log-likelihood function from the aspect model, and we need to maximize it.

$$\mathcal{L} = \sum_{d \in \mathcal{D}} \sum_{w \in \mathcal{W}} n(d, w) \log P(d, w),$$

- Expectation Maximization (EM) is used for this purpose

EM Steps

- E-Step
 - Expectation step where expectation of the likelihood function is calculated with the current parameter values
- M-Step
 - Update the parameters with the calculated posterior probabilities
 - Find the parameters that maximizes the likelihood function

E Step

- It is the probability that a word w occurring in a document d , is explained by aspect z

$$P(z|d, w) = \frac{P(z)P(d|z)P(w|z)}{\sum_{z'} P(z')P(d|z')P(w|z')},$$

(based on some calculations)

M Step

- All these equations use $p(z|d,w)$ calculated in E Step

$$P(w|z) = \frac{\sum_d n(d, w)P(z|d, w)}{\sum_{d, w'} n(d, w')P(z|d, w')},$$

$$P(d|z) = \frac{\sum_w n(d, w)P(z|d, w)}{\sum_{d', w} n(d', w)P(z|d', w)},$$

$$P(z) = \frac{1}{R} \sum_{d, w} n(d, w)P(z|d, w), \quad R \equiv \sum_{d, w} n(d, w).$$

- Converges to local maximum of the likelihood function

“Arts”

“Budgets”

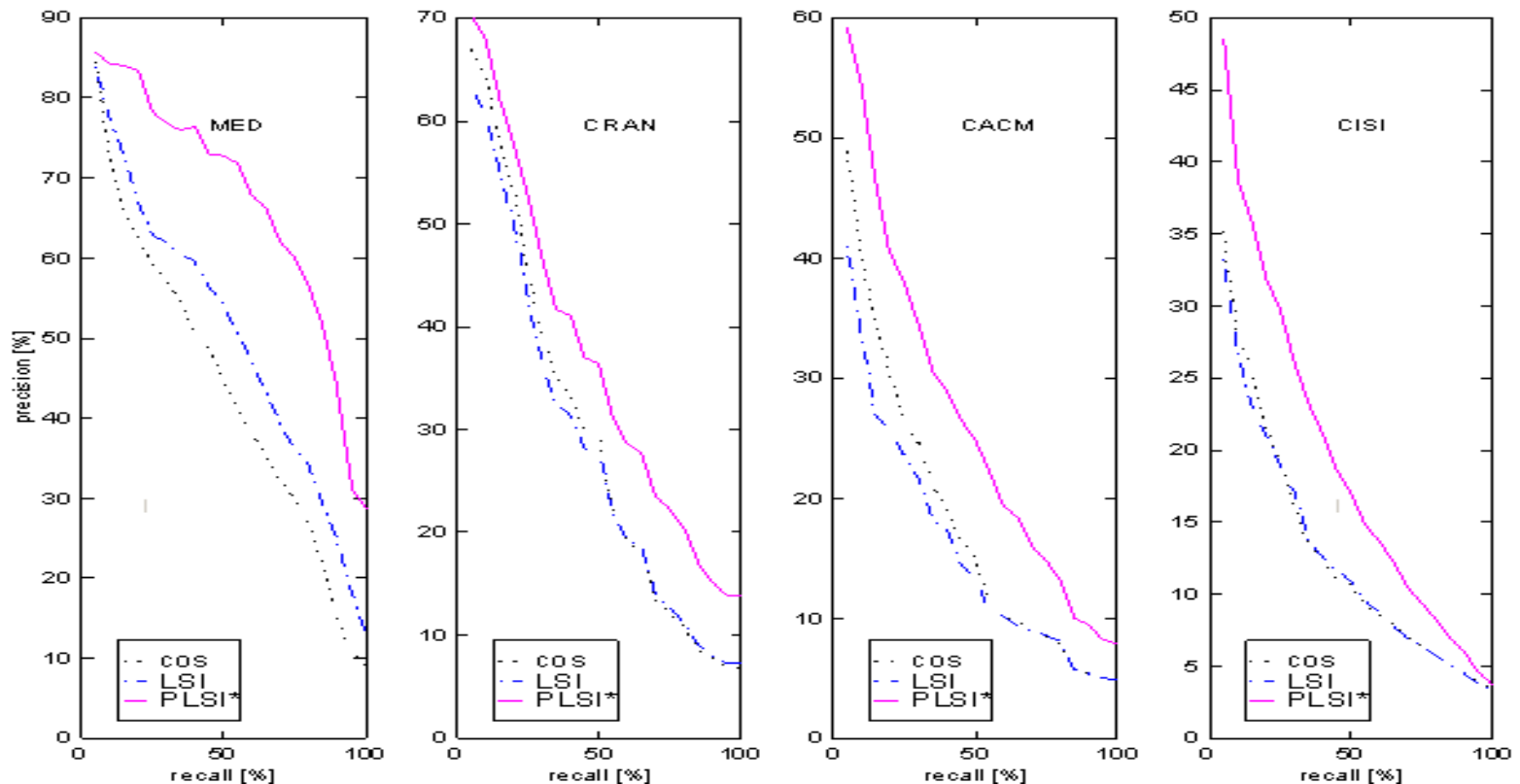
“Children”

“Education”

NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

The performance of a retrieval system based on this model (PLSI) was found superior to that of both the vector space based similarity (cos) and a non-probabilistic latent semantic indexing (LSI) method. (We skip details here.)



From Th. Hofmann, 2000

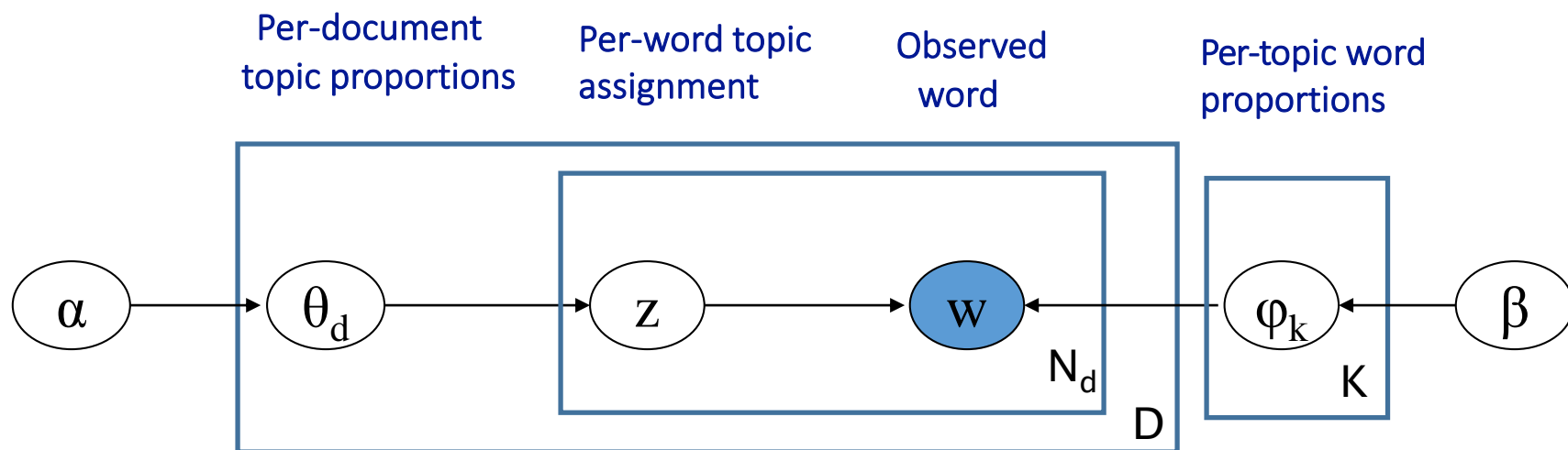
Comparing PLSA and LSA

- LSA and PLSA perform dimensionality reduction
 - In LSA, by keeping only K singular values
 - In PLSA, by having K aspects
- Comparison to SVD
 - U Matrix related to $P(d|z)$ (doc to aspect)
 - V Matrix related to $P(z|w)$ (aspect to term)
 - E Matrix related to $P(z)$ (aspect strength)
- The main difference is the way the approximation is done
 - PLSA generates a model (aspect model) and maximizes its predictive power
 - Selecting the proper value of K is heuristic in LSA
 - Model selection in statistics can determine optimal K in PLSA

Latent Dirichlet Allocation

• Generative Model

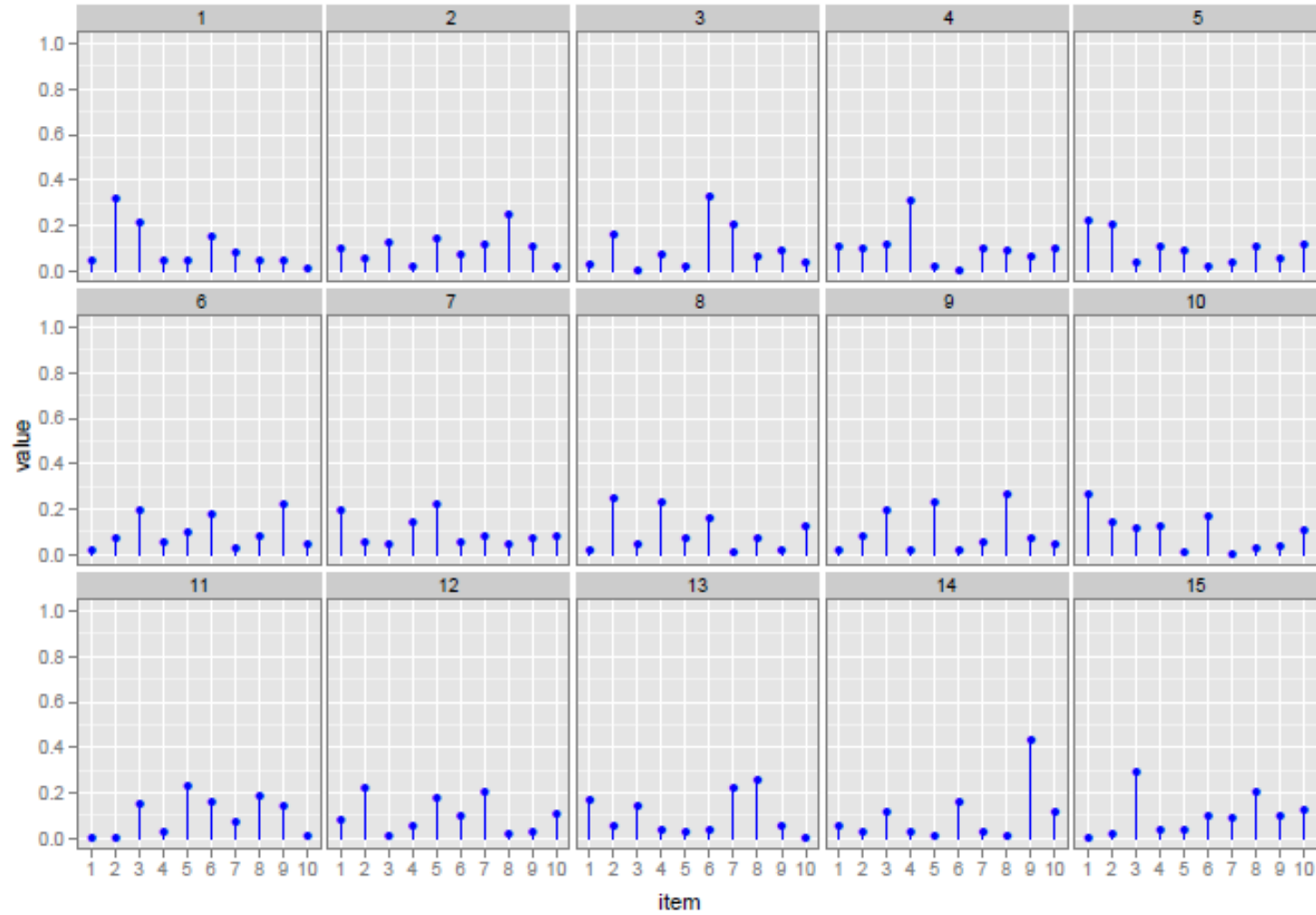
1. Choose $\theta_i \sim \text{Dir}(\alpha)$, where $i \in \{1, \dots, M\}$ and $\text{Dir}(\alpha)$ is a Dirichlet distribution
2. Choose $\varphi_k \sim \text{Dir}(\beta)$, where $k \in \{1, \dots, K\}$ and β typically is sparse
3. For each of the word positions i, j , where $j \in \{1, \dots, N_i\}$, and $i \in \{1, \dots, M\}$
 - (a) Choose a topic $z_{i,j} \sim \text{Multinomial}(\theta_i)$.
 - (b) Choose a word $w_{i,j} \sim \text{Multinomial}(\varphi_{z_{i,j}})$.



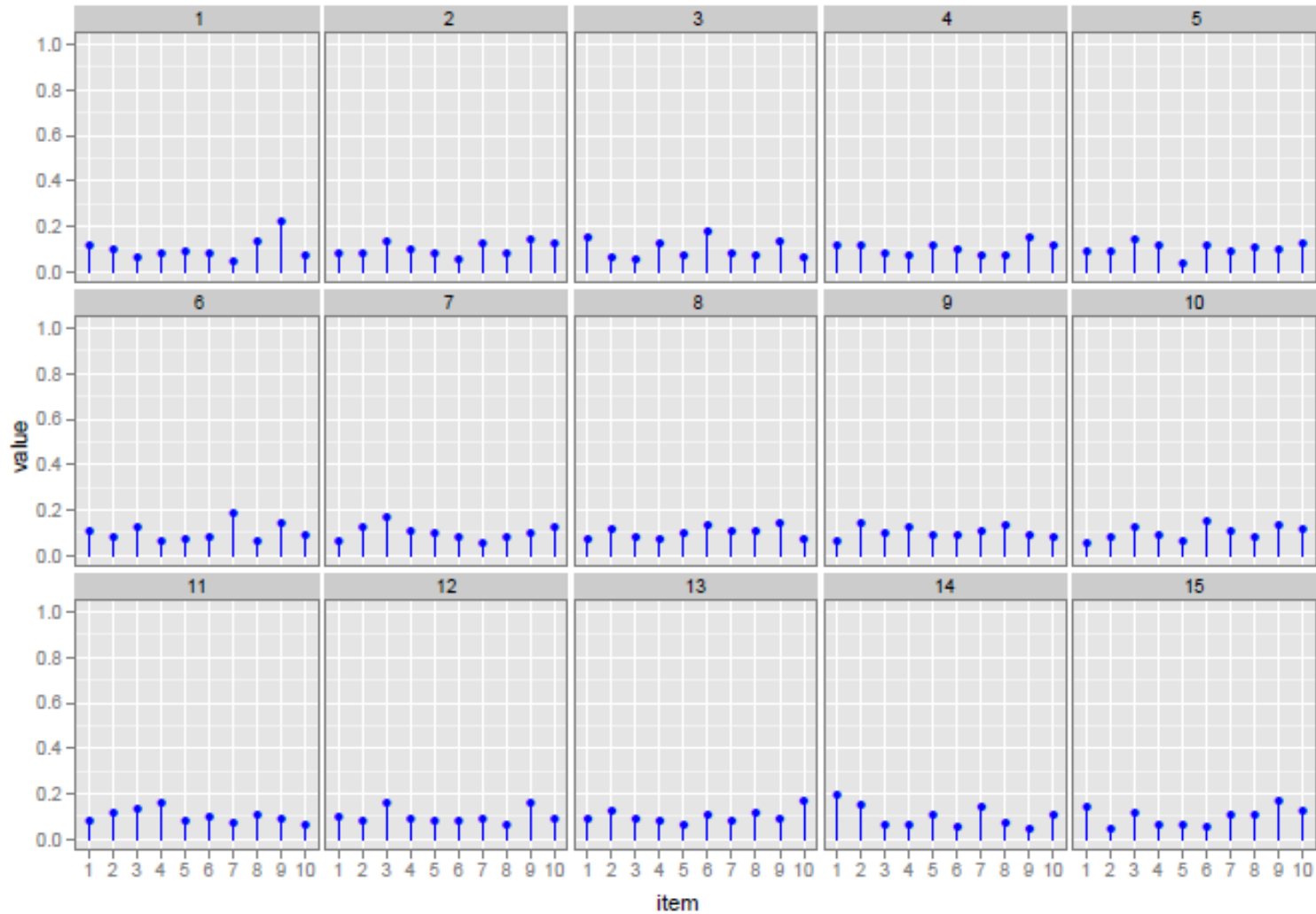
Joint Distribution

$$\prod_{i=1}^K p(\boldsymbol{\varphi}_i | \boldsymbol{\beta}) \prod_{d=1}^D p(\theta_d | \alpha) \left(\prod_{n=1}^N p(z_{d,n} | \theta_d) p(w_{d,n} | \boldsymbol{\varphi}_{1:K}, z_{d,n}) \right)$$

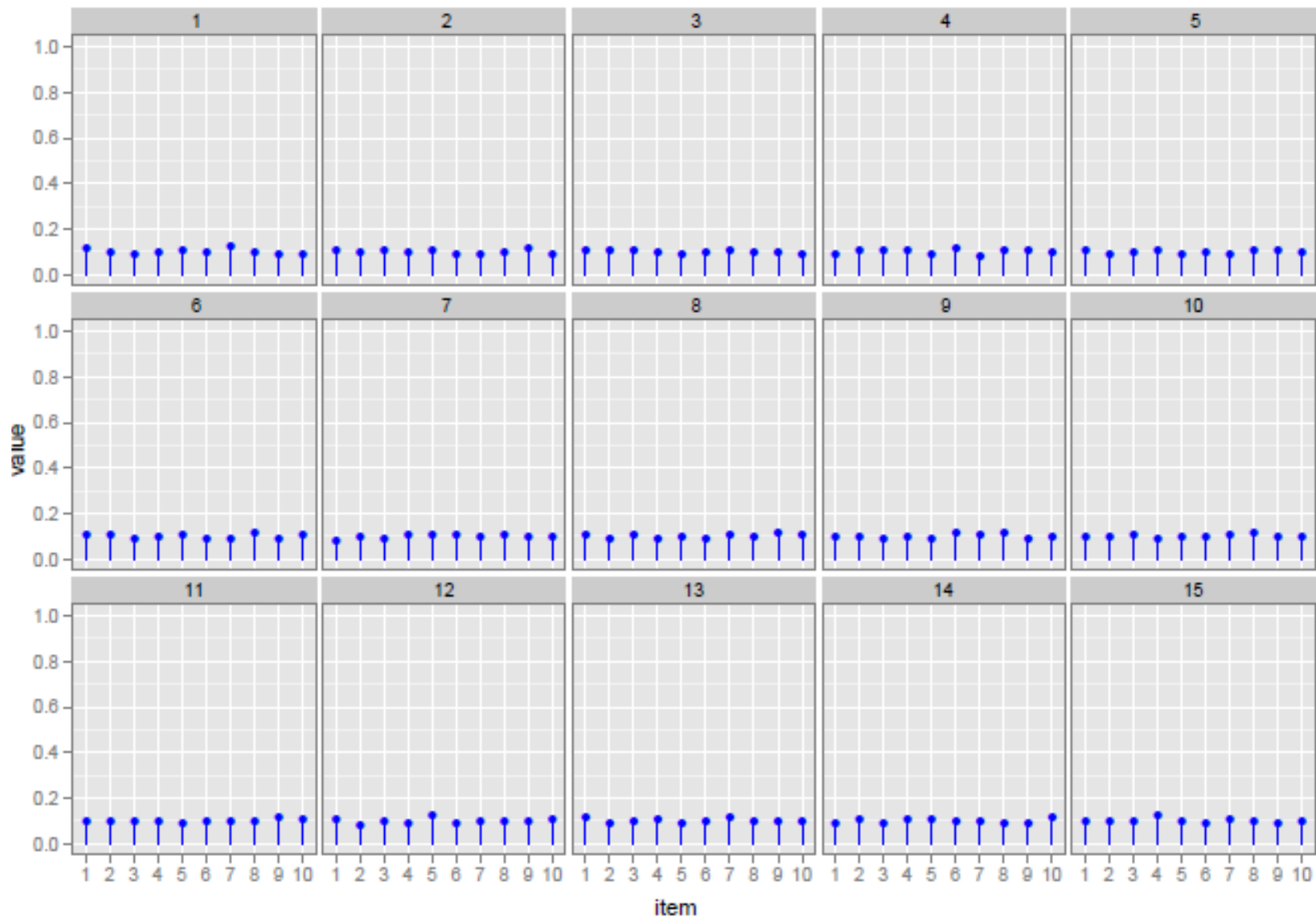
$$\alpha = 1$$



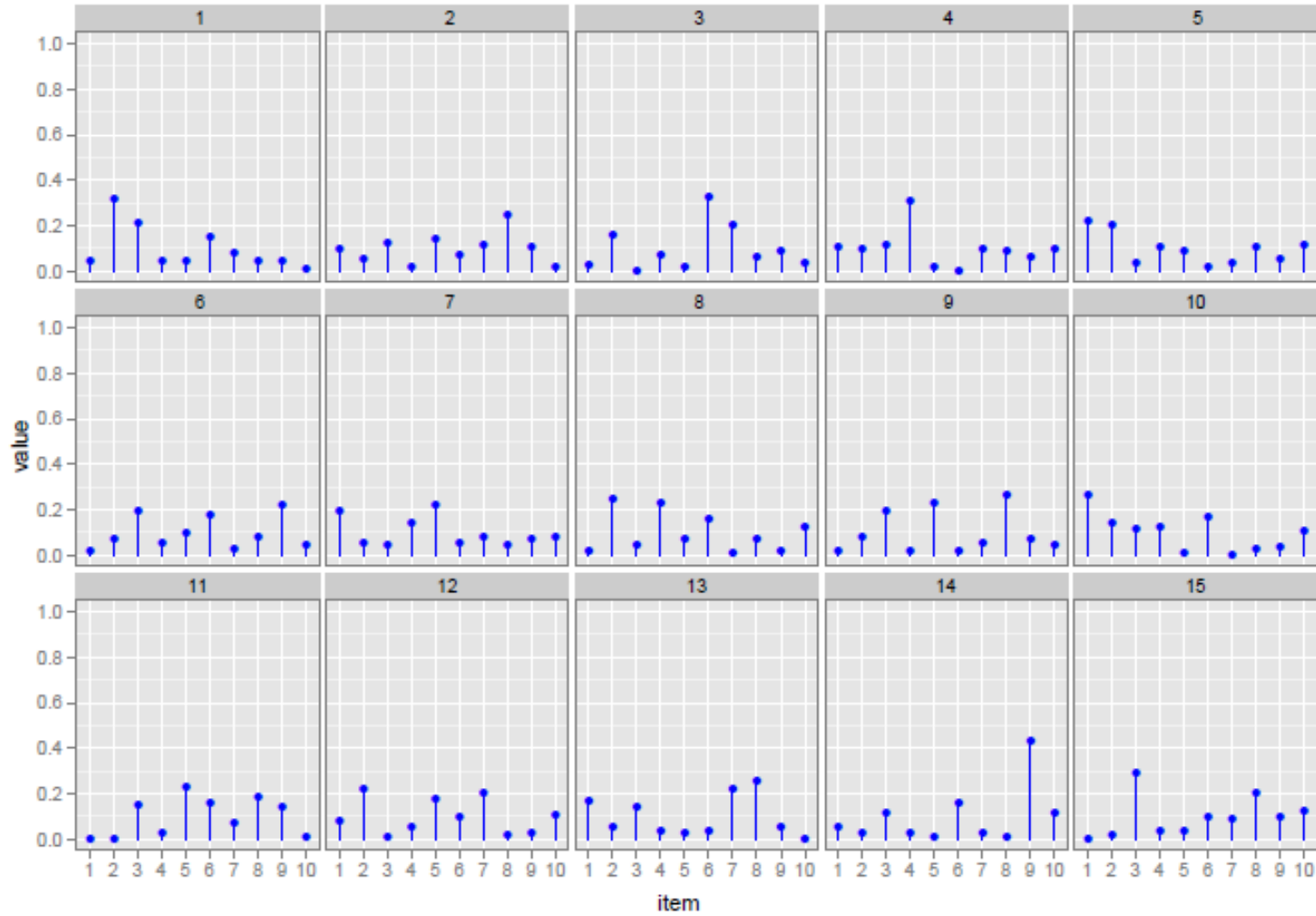
$$\alpha = 10$$



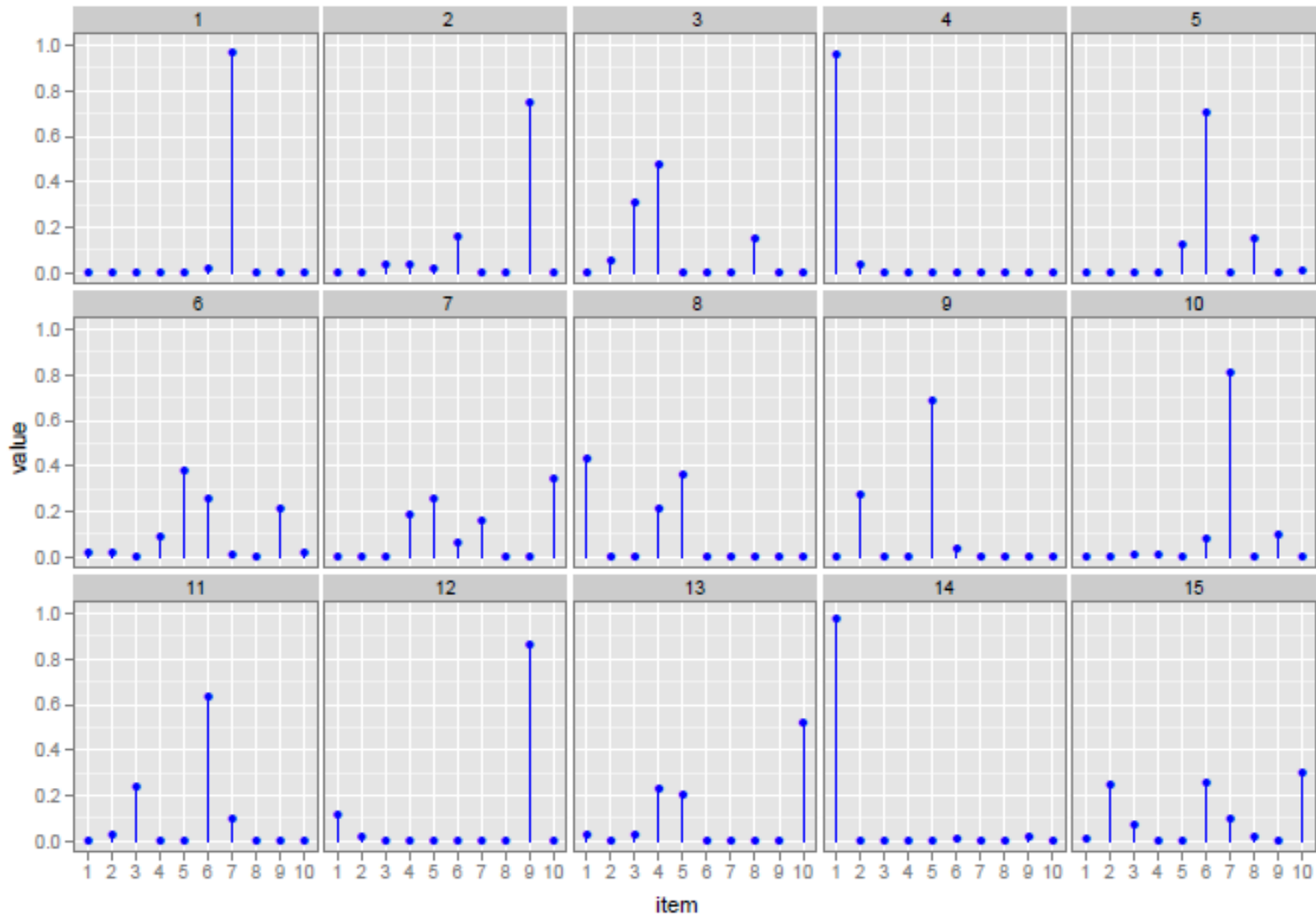
$$\alpha = 100$$



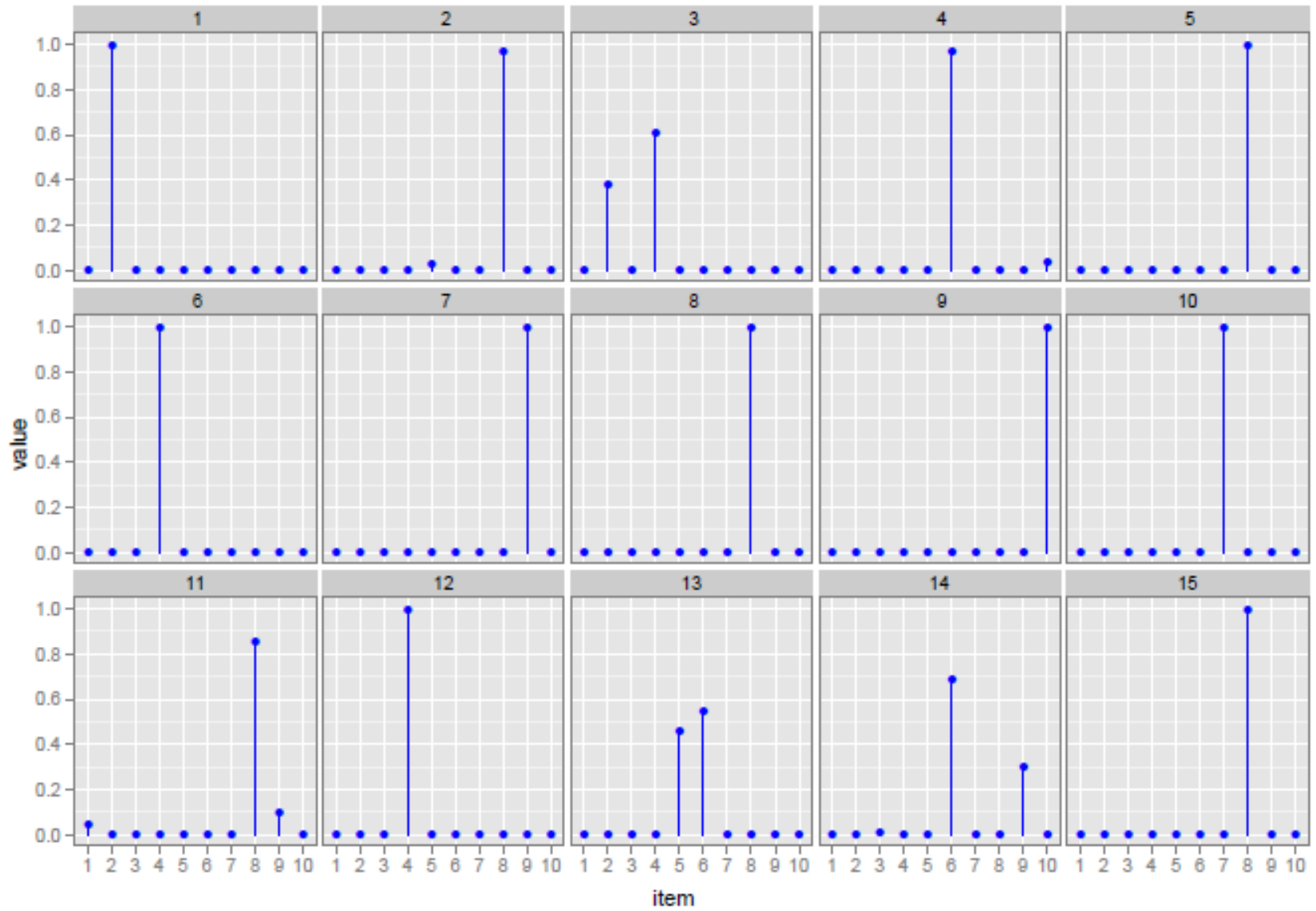
$$\alpha = 1$$



$$\alpha = 0.1$$



$$\alpha = 0.01$$



LDA vs. PLSA

- Explicit modeling of sparsity
- LDA was a rage
 - Lots of tools
- Often gives similar results