



# variable latent semantic indexing

ravi kumar

yahoo! research

ravikumar@yahoo-inc.com



joint work with

---

- anirban dasgupta, **cornell university**
- prabhakar raghavan, yahoo! research
- andrew tomkins, yahoo! research

kdd 2005, to appear



# outline

---

1. vector-space model
2. latent semantic indexing
3. variable latent semantic indexing

**sources: berry, hofmann, langville, raghavan**



# thesaurus approach

---

luhn (1957, 1961)

- similar/related words grouped into **notional families**
- represent document by notional elements (vocabulary based)
- matching by measuring degree of notional similarity
- key word in context indexing (**kwic**); technical manuals

quick <b>brown</b> fox	line 1
lazy <b>dog</b>	...
brown <b>fox</b> jumps	...
fox <b>jumps</b> over	...
the <b>lazy</b> dog	...



# lattices of search terms

---

joyce, needham (1958)

- lattices and hierarchies of search terms
  - thesaurus/dictionary
- if a term occurs in a document, then all all terms present at higher levels of the lattices are also considered occurring
  - punch cards!



# term associations

---

doyle (1962)

- identify unusual co-occurrences of pairs of words
  - term associations
- statistical methods used to identify correlations
  - chi-square tests, pearson correlation coefficient
- term association maps for interactive retrieval
- semantic maps
  - natural language understanding



# vector-space model

---

salton (1960s)

- each dictionary term is a dimension
- each document is a high-dimensional vector
- entries can be boolean or weighted
- geometric interpretation for similarity, etc
- **postulate:** documents close in vector space talk about same things
- can implement “query-by-example”



# term-document matrix

- $D = \{ d_1, \dots \}$  = set of **documents** in the collection
- $T = \{ t_1, \dots \}$  = set of **terms**

	$d_1$		$d_j$	
$t_1$				
$t_i$			$a_{ij}$	



# selecting terms

---

## pseudo-linguistic heuristics

- drop stop-words  
eg, articles, “html”, ...
- stemming/lemmatization  
eg, speaks, speaking  $\alpha$  speak
- use only nouns/noun phrases  
people, places, entity extraction
- phrase extraction  
eg, “university of california”
- domain-specific knowledge



# similarity measure

---

similarity between document and query

euclidean distance: sensitive to vector length

**angular cosine measure**

$$\text{sim}(d, q) = \cos \hat{A}(d, q) = \mathbf{hd}, \mathbf{q} / (\|\mathbf{kd}\| \text{ } \& \mathbf{kqk})$$

given a query

compute  $\delta_i = \text{sim}(d_i, q)$

rank documents according to  $\delta_i$

retrieve top documents



# example from berry's book

---

## terms

t<sub>1</sub>: bab(y,ies,y's)

t<sub>2</sub>: child(ren's)

t<sub>3</sub>: guide

t<sub>4</sub>: health

t<sub>5</sub>: home

t<sub>6</sub>: infant

t<sub>7</sub>: proofing

t<sub>8</sub>: safety

t<sub>9</sub>: toddler

## documents

d<sub>1</sub>: **infant & toddler** first aid

d<sub>2</sub>: **babies** and **children's** room  
(for your **home**)

d<sub>3</sub>: **child safety** at home

d<sub>4</sub>: your **baby's health & safety:**  
from **infant** to **toddler**

d<sub>5</sub>: **baby proofing** basics

d<sub>6</sub>: your **guide** to easy rust  
**proofing**

d<sub>7</sub>: beanie **babies** collector **guide**



eg, query = baby health

$$\mathbf{A} = \begin{matrix} & d_1 & d_2 & d_3 & d_4 & d_5 & d_6 & d_7 \\ \begin{matrix} t_1 \\ t_2 \\ t_3 \\ t_4 \\ t_5 \\ t_6 \\ t_7 \\ t_8 \\ t_9 \end{matrix} & \begin{pmatrix} 0 & 1 & 0 & 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix} \quad \mathbf{q} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \delta = \begin{bmatrix} \delta_1 \\ \delta_2 \\ \delta_3 \\ \delta_4 \\ \delta_5 \\ \delta_6 \\ \delta_7 \end{bmatrix} = \begin{bmatrix} 0 \\ .5774 \\ 0 \\ .8944 \\ .7071 \\ 0 \\ .7071 \end{bmatrix}$$

**answer** = d4: your baby's health & safety: from infant to toddler



**smart system**

---

**system for mechanical analysis and retrieval of text**

**aka**

**salton's magical automatic retriever of text**



# term-weighting

---

- document length normalization
- global term weights
  - boolean
  - tf-idf
  - okapi



# advantages of vector-space model

---

- no subjective selection of index terms
- partial match
  - no document contains all search terms
- highly automatable
  - similarity scores
- flexible term-weighting schemes
- extensions
  - document clustering (doc-doc similarity)
  - term clustering
  - relevance feedback
- geometric foundation



# synonymy

---

**synonymy**: different terms with similar meaning

- car vs automobile
- words indicating the same topic
- vector space model does not identify them together
- context can be used to bring them together

... car ...

... road ...

... accident ...

... hurt ...

... gas ...

... road ...

... automobile ...

... gas ...

... injury ...

... accident ...



# polysemy

---

**polysemy**: same word with different meanings  
depending on context

- jaguar, windows, apple
- esp severe in heterogeneous collection of documents
- vector space model does not discriminate them
- context can be used to differentiate them

... hunts ...

... jaguar ...

... wild ...

... nocturnal ...

... expensive ...

... mph ...

... jaguar ...

... car ...



# drawbacks of vector space model

---

- **dimensionality**
  - representation of documents is in very high dimension
  - curse of dimensionality for learning/estimation
- **document vectors**
  - sparse
  - cosine similarity could be noisy
- **syntax**
  - bag-of-words
  - phrases/order-relationships missed out
- **semantics**
  - only exact matches identified
  - semantic relations between words not exploited



# latent semantic indexing

---



# latent semantic indexing

---

- general idea
  - map documents (and terms) to a low-dimensional representation
  - design a mapping such that the low-dimensional space reflects semantic associations (**latent semantic space**)
  - compute document similarity based on the inner product in this latent semantic space
- goals
  - similar terms map to similar location in low dimensional space
  - noise reduction by dimension reduction



# dim. reduction vs. clustering

---

- **common use of dimension reduction**
  - find “better” representation of data
    - supporting more accurate retrieval
    - supporting more efficient retrieval
  - still using all points, but in a new representational space
- **common use of clustering**
  - summarize data or reduce data to fewer objects
  - clusters are often first-class citizens, directly used in the user interface or as part of retrieval algorithm



# random projections

---

take vectors and pack them into small number of dimensions

- choose a random direction  $\mathbf{x}_1$  in the vector space
- for  $i = 2$  to  $k$ 
  - choose random direction  $\mathbf{x}_i$  ?  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{i-1}$
- project each document vector into the subspace spanned by  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$

**johnson-lindenstrauss theorem:** with high probability, relative distances are (approximately) preserved by projection



# linear algebra recap

---



# eigenvalues and eigenvectors

**eigenvectors**, for a square  $m \times m$  matrix  $S$

$$S v = \lambda v$$

(right) eigenvector  $v \in \mathbb{R}^m \neq 0$       eigenvalue  $\lambda \in \mathbb{R}$

example

$$\begin{pmatrix} 6 & -2 \\ 4 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \end{pmatrix} = 2 \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$$Sv = \lambda v, (S - \lambda I)v = 0$$

only has a non-zero solution if  $|S - \lambda I| = 0$

this is an  $m$ -th order equation in  $\lambda$  that can have at most  $m$  distinct solutions (roots of the characteristic polynomial) – can be complex even though  $S$  is real



# matrix-vector multiplication

---

$$S = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

has eigenvalues 3, 2, 0 with corresponding eigenvectors

$$v_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad v_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad v_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

on each eigenvector,  $S$  acts as a multiple of the identity matrix: but as a different multiple on each.

any vector can be viewed as a linear combination of the eigenvectors eg,  $(2 \ 4 \ 6)^t = 2v_1 + 4v_2 + 6v_3$



# matrix-vector multiplication

---

- thus a matrix-vector multiplication such as  $Sx$  can be rewritten in terms of the eigenvalues/vectors

$$Sx = S(2v_1 + 4v_2 + 6v_3)$$

$$Sx = 2Sv_1 + 4Sv_2 + 6Sv_3 = 2\lambda_1 v_1 + 4\lambda_2 v_2 + 6\lambda_3 v_3$$

- even though  $x$  is an arbitrary vector, the action of  $S$  on  $x$  is determined by the eigenvalues/vectors
- **suggestion:** the effect of “small” eigenvalues is small



# facts: eigenvalues, eigenvectors

---

for symmetric matrices, eigenvectors for distinct eigenvalues are orthogonal

$$Sv_1 = \lambda_1 v_1, Sv_2 = \lambda_2 v_2, \lambda_1 \neq \lambda_2 \Rightarrow \langle v_1, v_2 \rangle = 0$$

all eigenvalues of a real symmetric matrix are real

all eigenvalues of a positive semi-definite matrix are non-negative

$$\forall w \in \mathbb{R}^n, w^t S w \geq 0, \text{ then if } Sv = \lambda v \Rightarrow \lambda \geq 0$$



## example

---

- let  $S = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$  ← real, symmetric

- the eigenvalues are 1 and 3 (nonnegative, real)
- the eigenvectors are orthogonal (and real)

$$\begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$



# eigen/diagonal decomposition

---

let  $S \in \mathbb{R}^{m \times m}$  be a square matrix with  $m$  linearly independent eigenvectors (a “non-defective” matrix)

- **matrix diagonalization theorem:** exists an eigen decomposition

$$S = U \Lambda U^{-1}$$

- columns of  $U$  are eigenvectors of  $S$
- diagonal elements of  $\Lambda$  are eigenvalues of  $S$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n), \lambda_i \neq \lambda_{i+1}$$



## proof of theorem

---

let  $U$  have the eigenvectors as columns  $U = [v_1, \dots, v_n]$

then,  $SU$  can be written

$$\begin{aligned} S U &= S [v_1, \dots, v_n] = [\lambda_1 v_1, \dots, \lambda_n v_n] \\ &= [v_1, \dots, v_n] \Phi \text{diag}(\lambda_1, \dots, \lambda_n) \end{aligned}$$

thus

$$SU = U\Lambda$$

or

$$S = U \Lambda U^{-1}$$



## example

---

recall  $S = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}; \lambda_1 = 1, \lambda_2 = 3. \quad U = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$

inverting, we have  $U^{-1} = \begin{bmatrix} 1/2 & -1/2 \\ 1/2 & 1/2 \end{bmatrix}$

then,  $S = U\Lambda U^{-1} = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} 1/2 & -1/2 \\ 1/2 & 1/2 \end{bmatrix}$



## eg, contd

---

divide  $U$  (and multiply  $U^{-1}$ ) by  $\sqrt{2}$

$$\text{then, } \mathbf{S} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

$\mathbf{Q} \qquad \qquad \Lambda \qquad \qquad (\mathbf{Q}^{-1} = \mathbf{Q}^t)$



# symmetric eigen decomposition

---

if  $S \in \mathbb{R}^{m \times m}$  is a symmetric matrix

**theorem:** exists a (unique) **eigen decomposition**

$$S = Q \Lambda Q^t$$

where  $Q$  is orthogonal

- $Q^{-1} = Q^t$
- columns of  $Q$  are normalized eigenvectors
- columns are orthogonal
- (everything is real)



# singular value decomposition

for an  $m \times n$  matrix  $A$  of rank  $r$  there exists a unique factorization (**singular value decomposition = svd**)

$$A = U \Sigma V^t$$

The diagram shows the equation  $A = U \Sigma V^t$  at the top. Below it, three boxes are arranged horizontally. The first box contains  $m \times n$ , the second contains  $m \times n$ , and the third contains  $V \text{ is } n \times n$ . Arrows point from each box to the corresponding matrix in the equation above: from the first box to  $U$ , from the second box to  $\Sigma$ , and from the third box to  $V^t$ .

columns of  $U$  are orthogonal eigenvectors of  $AA^t$

columns of  $V$  are orthogonal eigenvectors of  $A^tA$

eigenvalues  $\lambda_1 \dots \lambda_r$  of  $AA^T$  are the eigenvalues of  $A^T A$

$$\sigma_i = \sqrt{\lambda_i}$$

$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$  (singular values)



# singular value decomposition

illustration of svd dimensions and sparseness

$$\begin{aligned} \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_A &= \underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_U \underbrace{\begin{bmatrix} \bullet & & & & \\ & \bullet & & & \\ & & \bullet & & \\ & & & \bullet & \\ & & & & \bullet \end{bmatrix}}_\Sigma \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_{V^T} \\ \underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_A &= \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_U \underbrace{\begin{bmatrix} \bullet & & & & \\ & \bullet & & & \\ & & \bullet & & \\ & & & \bullet & \\ & & & & \bullet \end{bmatrix}}_\Sigma \underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_{V^T} \end{aligned}$$



## svd example

---

$$\text{let } A = \begin{bmatrix} 1 & -1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$$

thus  $m=3$ ,  $n=2$  and svd is

$$\begin{bmatrix} 0 & 2/\sqrt{6} & 1/\sqrt{3} \\ 1/\sqrt{2} & -1/\sqrt{6} & 1/\sqrt{3} \\ 1/\sqrt{2} & 1/\sqrt{6} & -1/\sqrt{3} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \sqrt{3} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix}$$

typically, the singular values arranged in decreasing order



# low-rank approximation

---

svd can be used to compute optimal low-rank approximations

**approximation problem:** find  $A_k$  of rank  $k$  such that

$$A_k = \min_{X: \text{rank}(X) = k} \|A - X\|_f$$

**frobenius norm:**  $\|A\|_f^2 = \sum_{i,j} |a_{ij}|^2$

$A_k$  and  $X$  are both  $m \times n$  matrices

typically, want  $k \ll r$

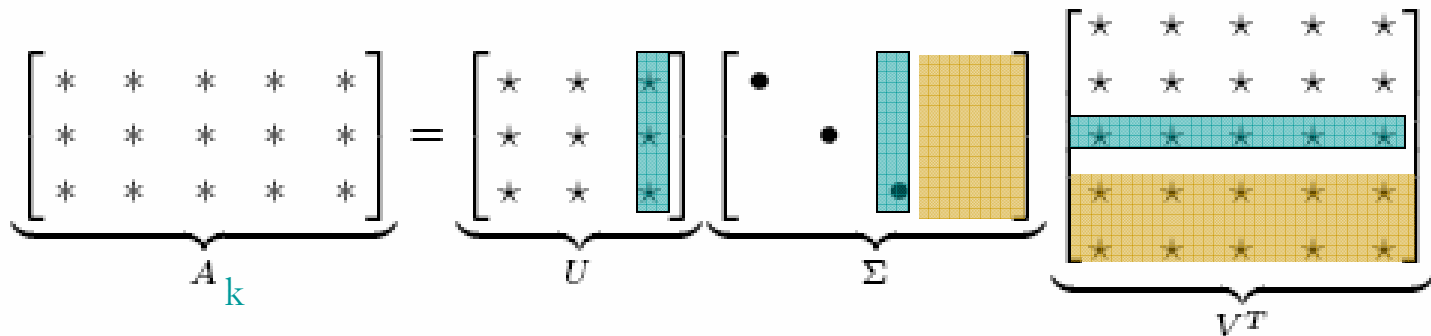


# low-rank approximation

solution via svd

$$A_k = U \underbrace{\text{diag}(\sigma_1, \dots, \sigma_k, 0, \dots, 0)} V^t$$

set smallest r-k  
singular values to zero



$$A_k = U_k \Sigma_k V_k^t = \sum_{i=1, k} \sigma_i u_i v_i^t = \text{sum of rank 1 matrices}$$



# approximation error

---

how good (bad) is this approximation?

it is the best possible, measured by the frobenius norm of the error

$$\min_{\mathbf{X}: \text{rank}(\mathbf{X}) = k} \|\mathbf{A} - \mathbf{X}\|_f = \|\mathbf{A} - \mathbf{A}_k\|_f = \sqrt{\sum_{j > k} \sigma_j^2}$$

where the  $\sigma_i$  are ordered such that  $\sigma_i \geq \sigma_{i+1}$

suggests why frobenius error drops as  $k$  increases

**eckart-young** theorem



# svd vs random projections

---

- completely different method for low-rank approximation
- random projection is data-oblivious
  - svd-based approximation is data-dependent
- error for random projection depended only on start/finish dimensionality
  - for every distance
- error for svd-based approximation is for the frobenius norm, not for individual distances
- eliminate redundant axes, brings relevant dimensions together



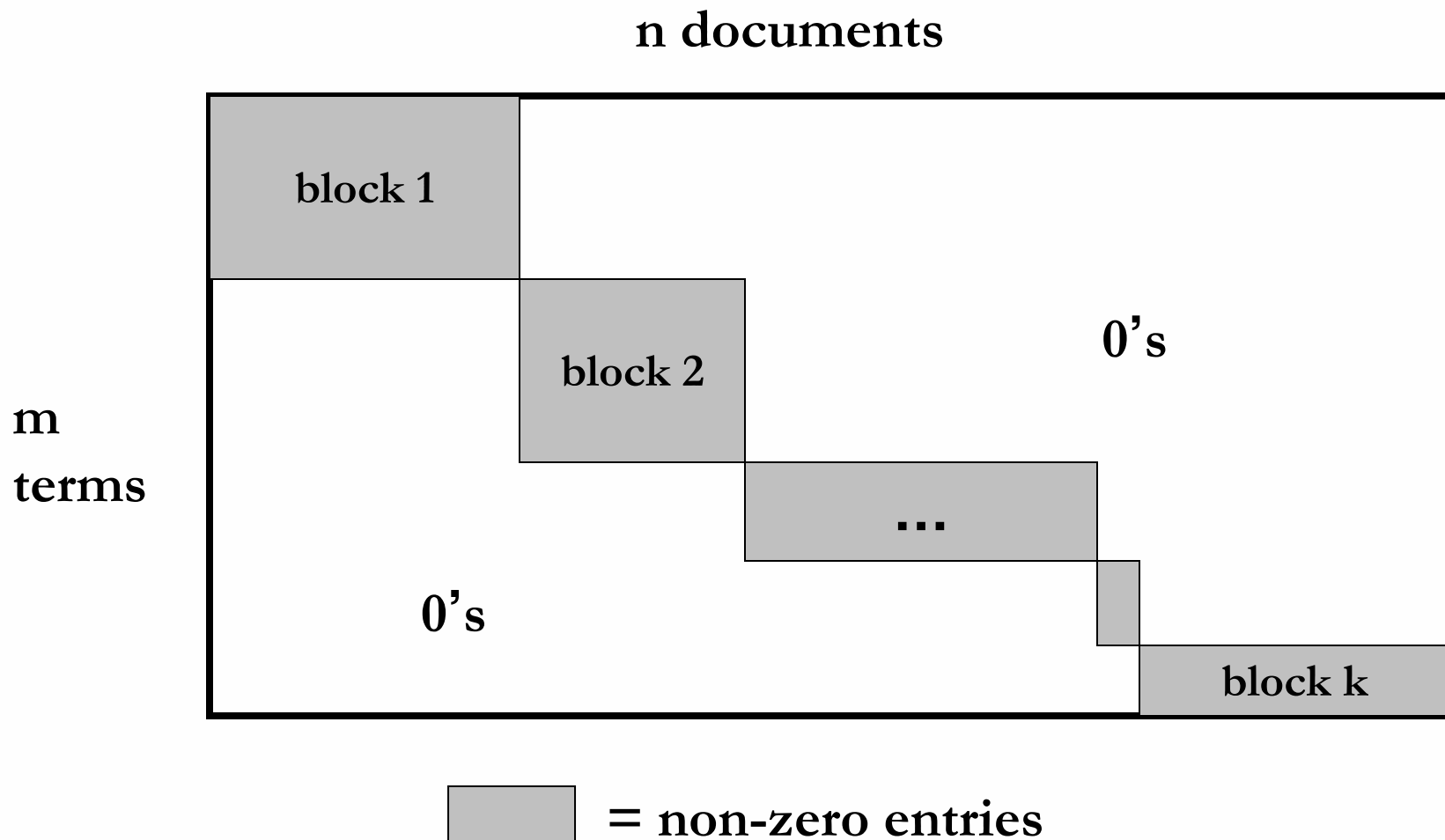
# latent semantic indexing (lsi)

---

- from term-document matrix  $A$ , compute the approximation  $A_k$
- there is a row for each term and a column for each document in  $A_k$
- thus documents live in a space of  $k \ll r$  dimensions
  - these dimensions are not the original axes
- query mapped to this space
  - not a sparse vector
- compute document similarity based on the inner product in this latent semantic space

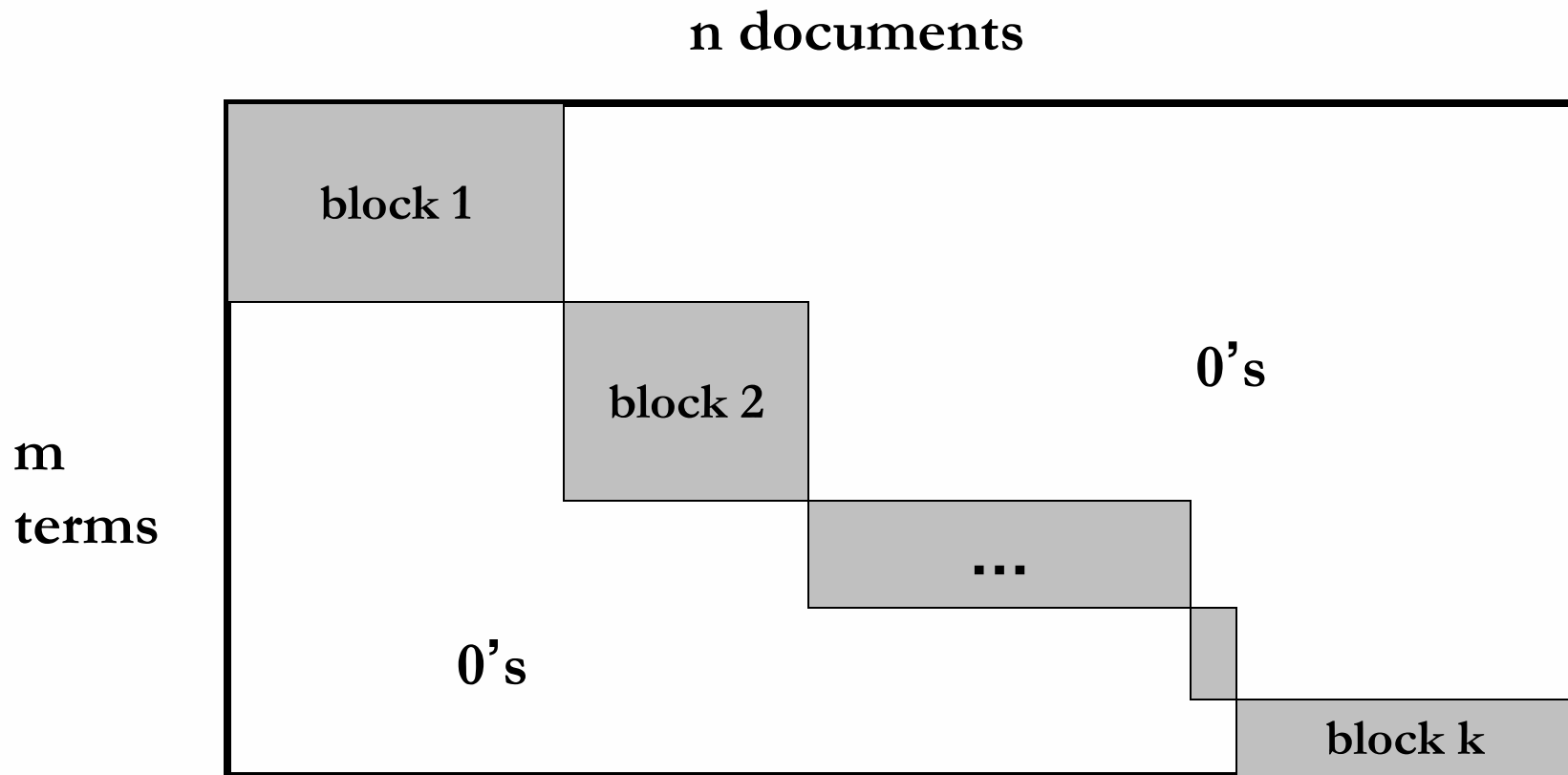


# intuition from block matrices





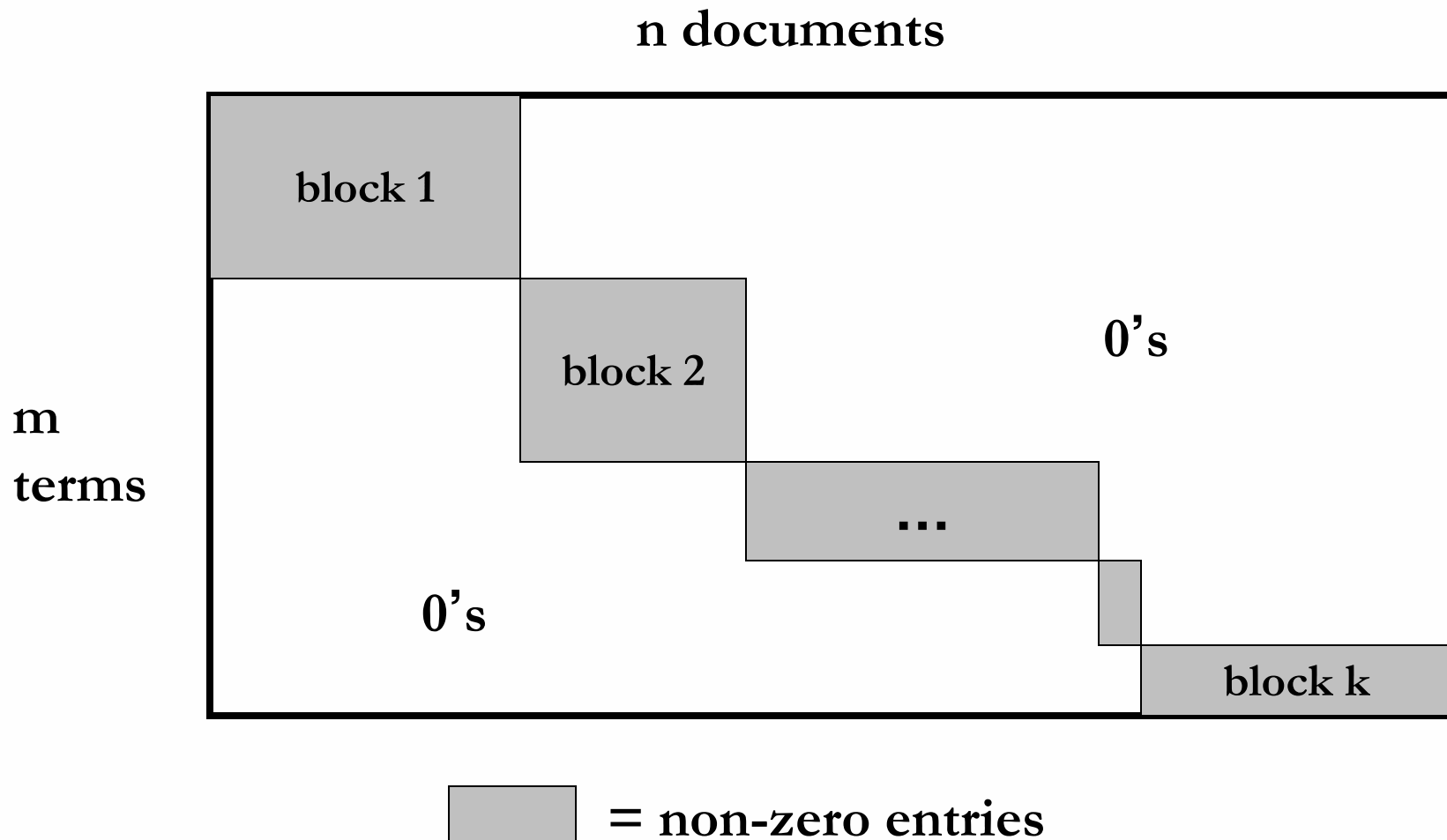
# intuition, contd



vocabulary partitioned into  $k$  topics (clusters); each document discusses only one topic

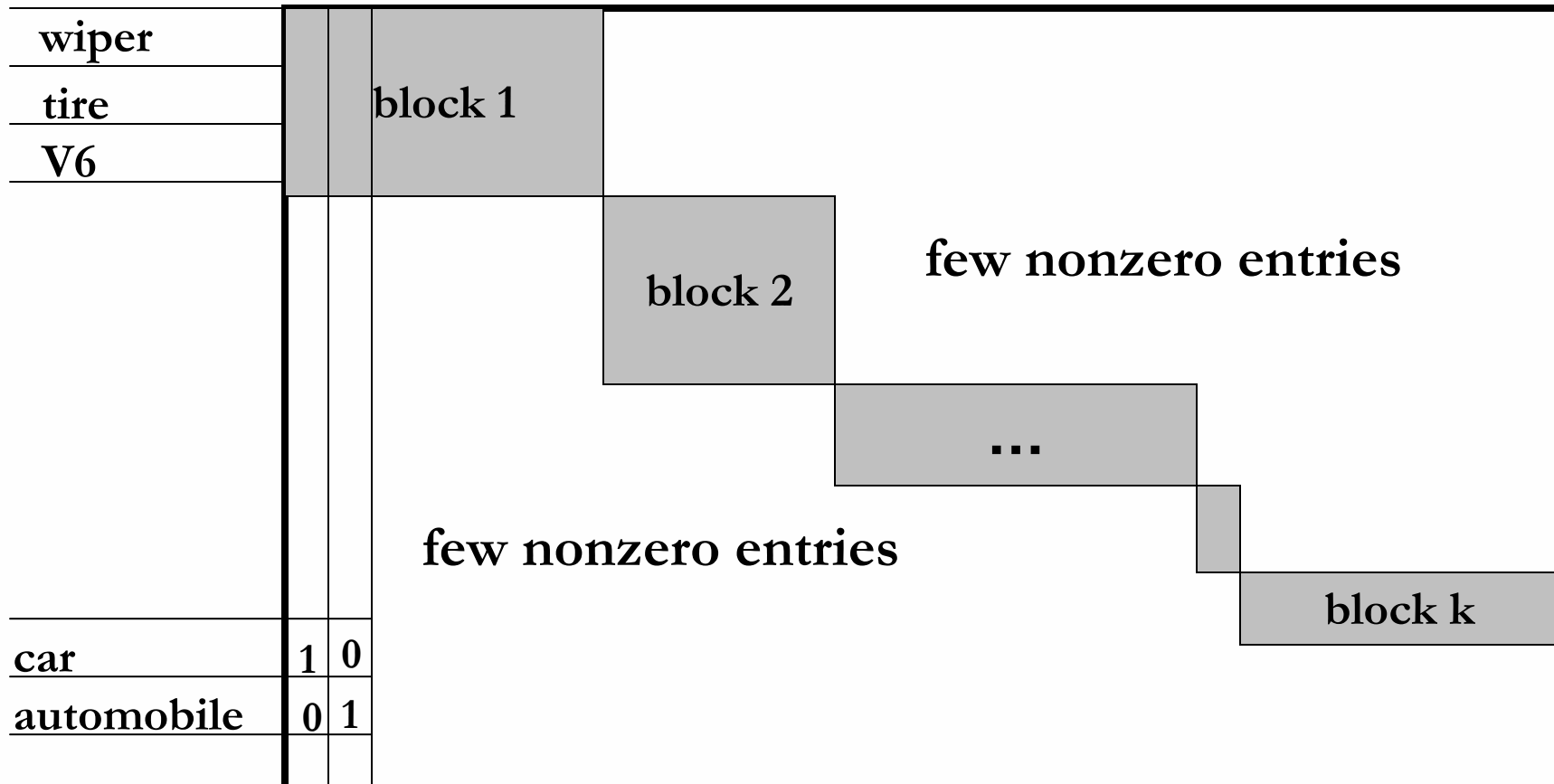


# intuition, contd





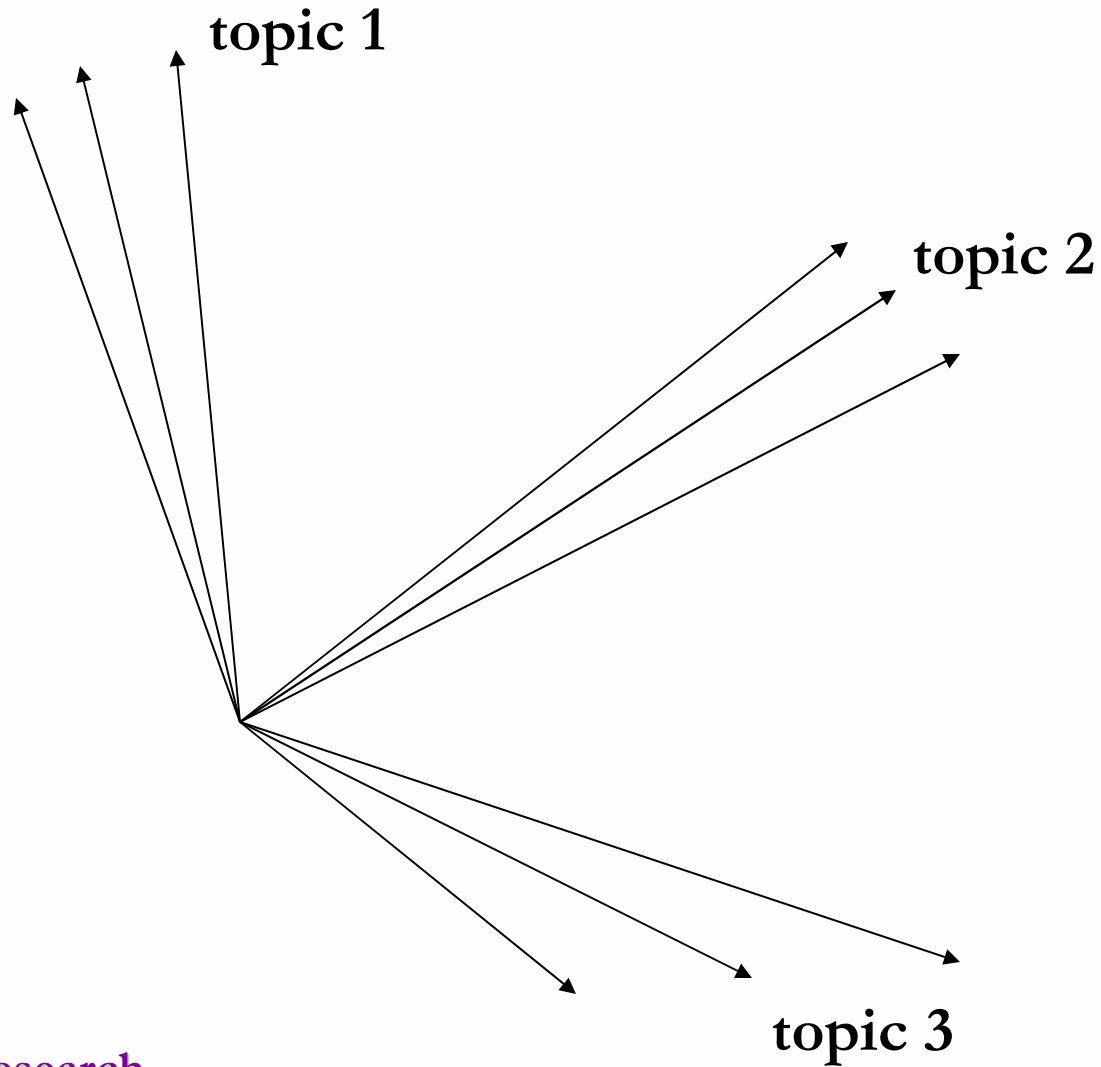
# intuition, contd





# simplistic picture

---





# some wild extrapolation

---

**“dimensionality” of a corpus is the number of distinct topics represented in it**

**more mathematical wild extrapolation**

- if  $A$  has a rank  $k$  approximation of low frobenius error, then there are no more than  $k$  distinct topics in the corpus**



# empirical evidence

---

- experiments on trec 1/2/3 – **dumais**
- lanczos svd code (available on netlib) due to **berry**  
used in these expts
  - running times of ~ one day on tens of thousands of docs
- dimensions – various values 250-350 reported
  - (under 200 reported unsatisfactory)
- generally expect recall to improve – what about precision?



## empirical evidence, contd

---

- precision at or above median trec precision
  - top scorer on almost 20% of trec topics
- slightly better on average than straight vector spaces
- effect of dimensionality

dimensions	precision
250	0.367
300	0.371
346	0.374



# variable latent semantic indexing

---



# low dimensional approx, recap

---

- extracts salient features from corpus
- low-dimensional representation
  - storage efficiency
  - computational efficiency
- disambiguation based on context
- “clustering” effect
  - brings relevant documents together
- noise reduction



## effect of queries

---

standard lsi is oblivious to queries

200-300 dimensions needed to get reasonable errors

suppose most of the queries come from finance, sports

can we tune lsi against this query distribution?

**ad hoc**: throw out non-finance, non-sport terms from  
dictionary and compute lsi

principled approach?



# linear algebra, recap

---

- svd

$$A = U \Sigma V^t$$

U, V column orthogonal,  $\Sigma$  diagonal

- symmetric psd eigen decomposition

$$C = Y \Lambda Y^t$$

Y orthogonal,  $\Lambda \succeq 0$

- frobenius norm

$$\|A\|_f^2 = \sum_{i,j} a_{ij}^2 = \text{Tr } A^t A$$

- svd low-rank approximation

$$A_k = \sum_{i=1, k} \sigma_i u_i v_i^t$$



# query distribution, co-occurrence

---

$Q$  = distribution on terms

**co-occurrence matrix**

$$C = C_Q = \mathbb{E}_{q \sim Q} [q q^t]$$

**fact:**  $C$  is positive semi-definite

**proof:**  $v^t C v = \mathbb{E}[v^t q q^t v] = \mathbb{E}[(v^t q)^2] \geq 0$

**corollary:**  $C = Y \Lambda Y^t$



# svd and average distortion

---

suppose  $Q$  is such that each coordinate is iid unit normal, ie,  $C = I$

find rank  $k$  matrix  $X$  that minimizes the average distortion imparted to random vector according to  $Q$

$$E_{q \sim Q} [ \| \mathbf{k} q^t (A - X) \mathbf{k}_2^2 ]$$

answer:  $X = A_k$



## proof when $C = I$

---

$$\begin{aligned} & \mathbf{E}[\mathbf{k} \mathbf{q}^t (\mathbf{A} - \mathbf{X}) \mathbf{k}_2^2] \\ &= \mathbf{E}[\mathbf{q}^t (\mathbf{A} - \mathbf{X})(\mathbf{A}^t - \mathbf{X}^t) \mathbf{q}] \\ &= \mathbf{E}[\text{Tr} ((\mathbf{A} - \mathbf{X})\mathbf{q}\mathbf{q}^t(\mathbf{A}^t - \mathbf{X}^t))] \\ &= \text{Tr} ((\mathbf{A}^t - \mathbf{X}^t) \mathbf{C} (\mathbf{A} - \mathbf{X})) \\ &= \text{Tr} ((\mathbf{A} - \mathbf{X})^t (\mathbf{A} - \mathbf{X})) \\ &= \mathbf{k}_{\mathbf{A} - \mathbf{X}}^2 \end{aligned}$$



## rotating the space

---

assume  $C$  is full rank and define

$$C^{1/2} = Y \Lambda^{1/2} Y^T$$

$$C^{-1/2} = Y \Lambda^{-1/2} Y^T$$

define  $w = C^{-1/2} q$

**fact:**  $E[ww^t] = I$

**proof:**  $E[ww^t] = E[C^{-1/2} qq^t C^{-1/2}]$   
 $= C^{-1/2} \mathbb{E} [C] C^{-1/2}$   
 $= I$



## main idea

---

- rotate the space of query vectors so that the query distribution becomes isotropic
- find the svd of term-document matrix in the rotated space
- compute the low-rank approximation
- rotate the low-rank approximation back to the original space



## some details

---

$$\begin{aligned} & \mathbf{E}[\mathbf{k} \mathbf{q}^t (\mathbf{A} - \mathbf{X}) \mathbf{k}^2] \\ &= \mathbf{E}[\mathbf{k} \mathbf{w}^t \mathbf{C}^{1/2} (\mathbf{A} - \mathbf{X}) \mathbf{k}^2] \\ &= \mathbf{E}[\mathbf{k} \mathbf{w}^t (\mathbf{C}^{1/2} \mathbf{A} - \mathbf{C}^{1/2} \mathbf{X}) \mathbf{k}^2] \end{aligned}$$

minimized when

$$\mathbf{C}^{1/2} \mathbf{X} = \text{svd} (\mathbf{C}^{1/2} \mathbf{A})$$

and so

$$\mathbf{X} = \mathbf{C}^{-1/2} \text{svd} (\mathbf{C}^{1/2} \mathbf{A})$$



## main result

---

### alternate statement

if  $V_k$  = top  $k$  right singular vectors of  $C^{1/2} A$  in its columns, then the minimizing solution is  $A V_k V_k^t$

space requirements = same as svd

store  $A V_k$  and  $V_k^t$



# experimental setup

---

- reuters data (1987)
  - 21k documents
  - five categories
  - 112k terms
  - 134 terms/document
- preprocessing
  - porter stemmed, case folded, stop worded
  - term-document matrix with boolean, okapi weighting
- svdpackc
  - external memory implementation



## metrics

---

- $l_2$  error
  - $\mathbf{k} \mathbf{q}^t (\mathbf{A} - \mathbf{A}_k) \mathbf{k}$
- **competitive precision**
  - $S =$  top  $k$  documents ranked by  $v_{ls}$
  - $cp(d) = |S \cap A[1\dots d]| / d$



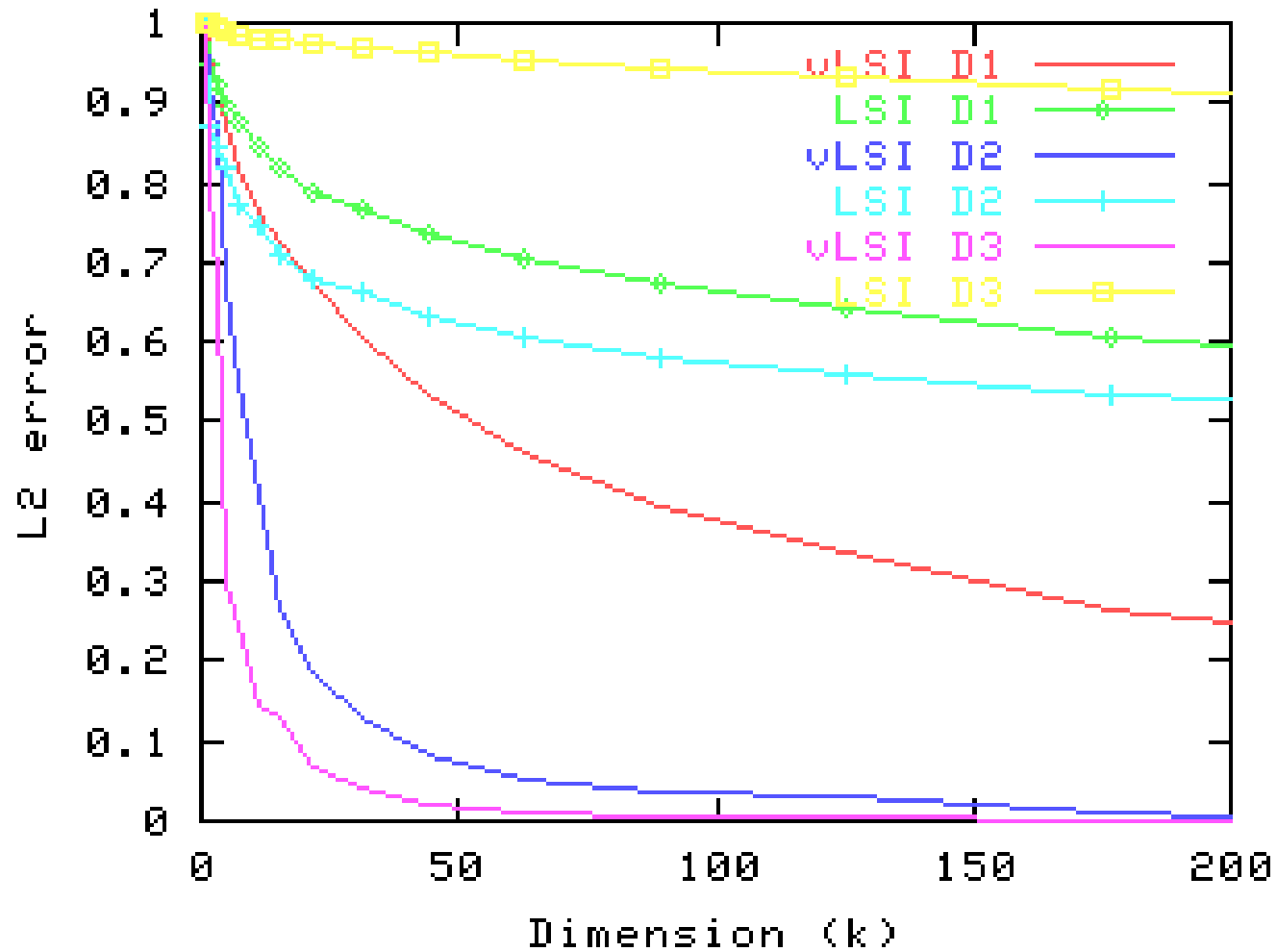
# query distributions

---

- single word
  - terms distributed according to distribution in corpus
  - terms distributed power law, ordered by distribution in corpus
  - terms distributed by power law, ordered randomly
- two topics
  - money, commodity
- double word
  - top two bigrams



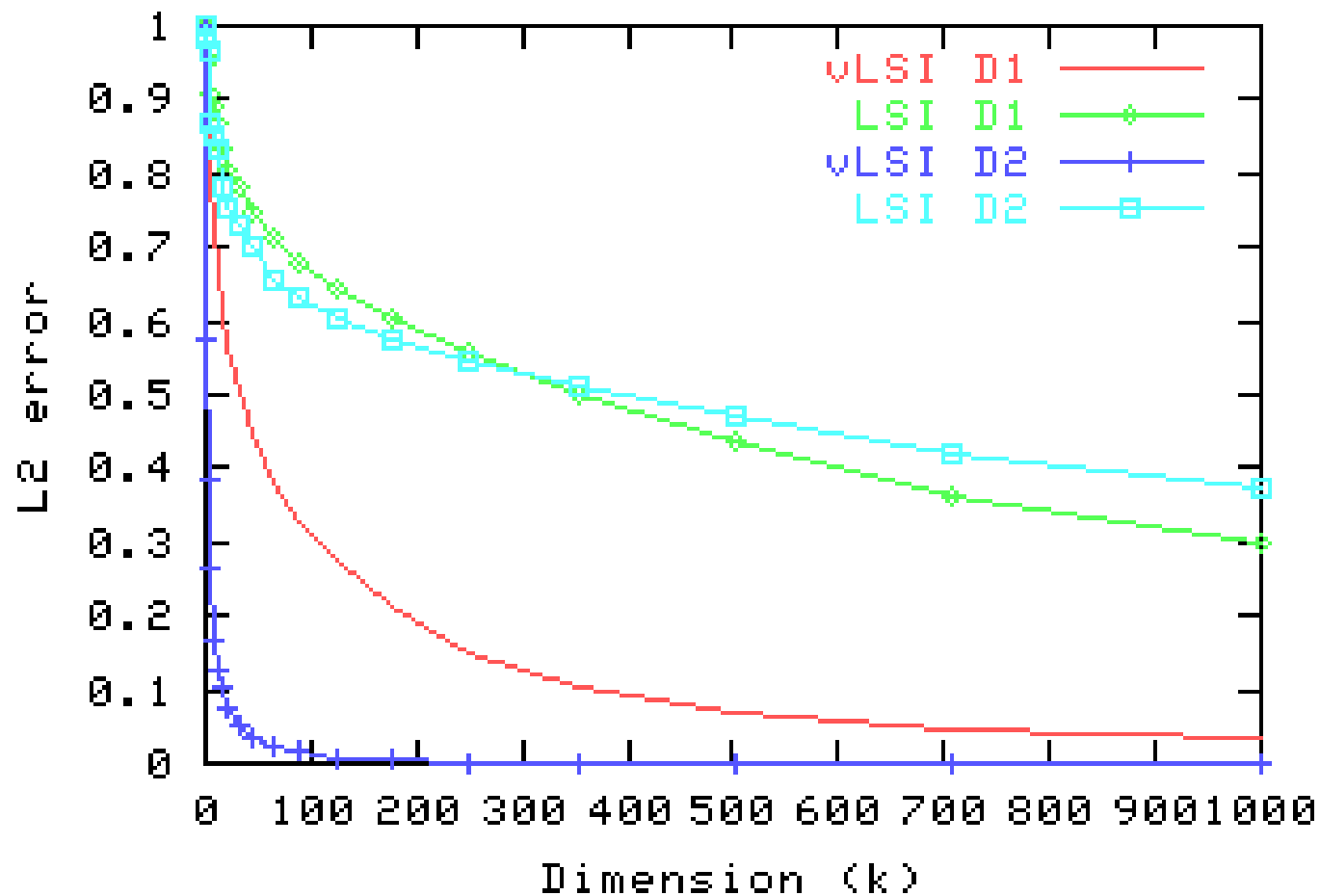
# results, l2 error





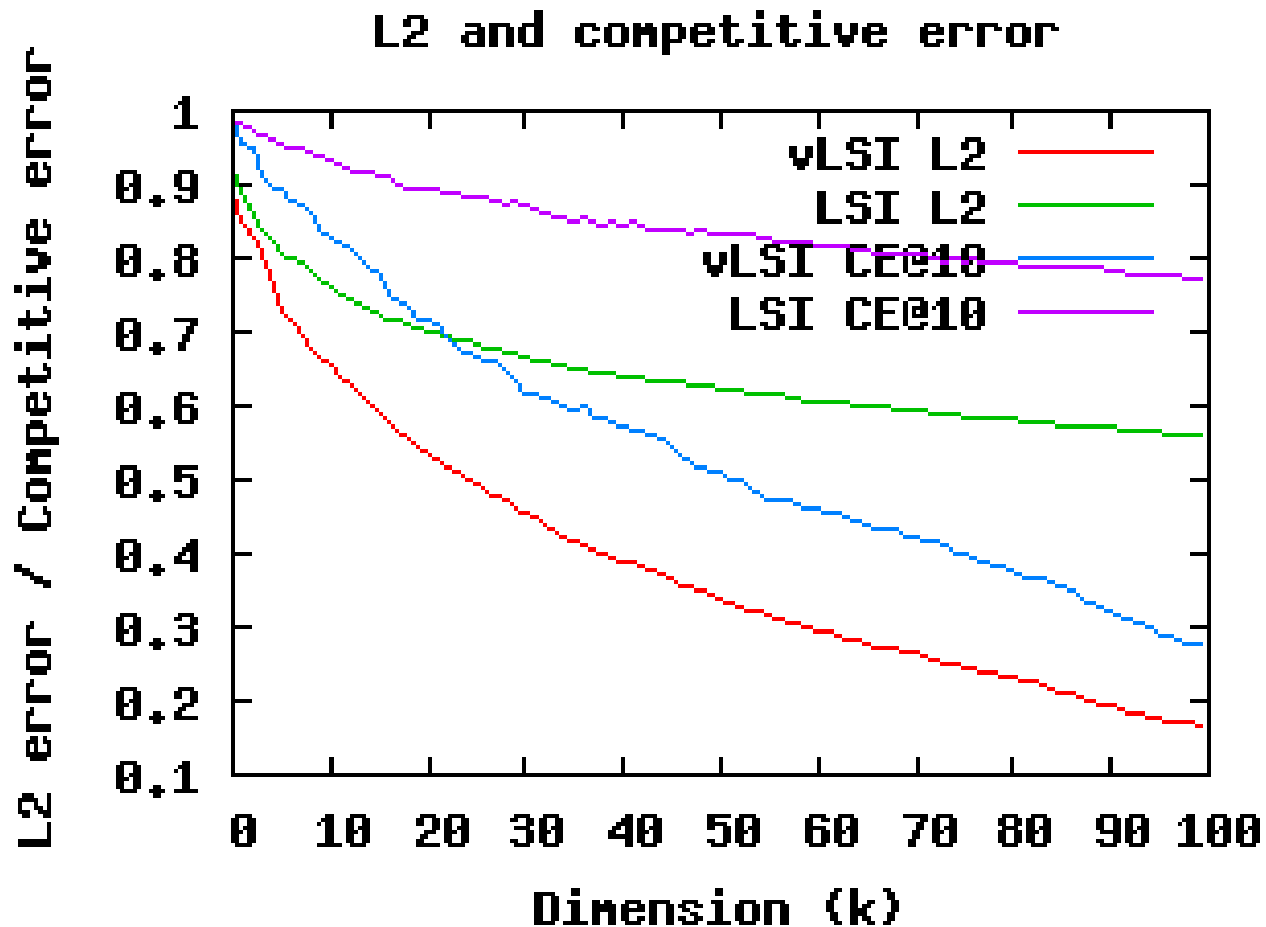
# results, money, l2 error

L2 error for money and commodities topics





# results, multiword queries





## future work

---

- personalized version
- computational questions
  - efficiency
  - approximate
  - sampling
  - data stream
- perturbation analysis
- beyond text domains



*thank you!*