

Dense Word Representation

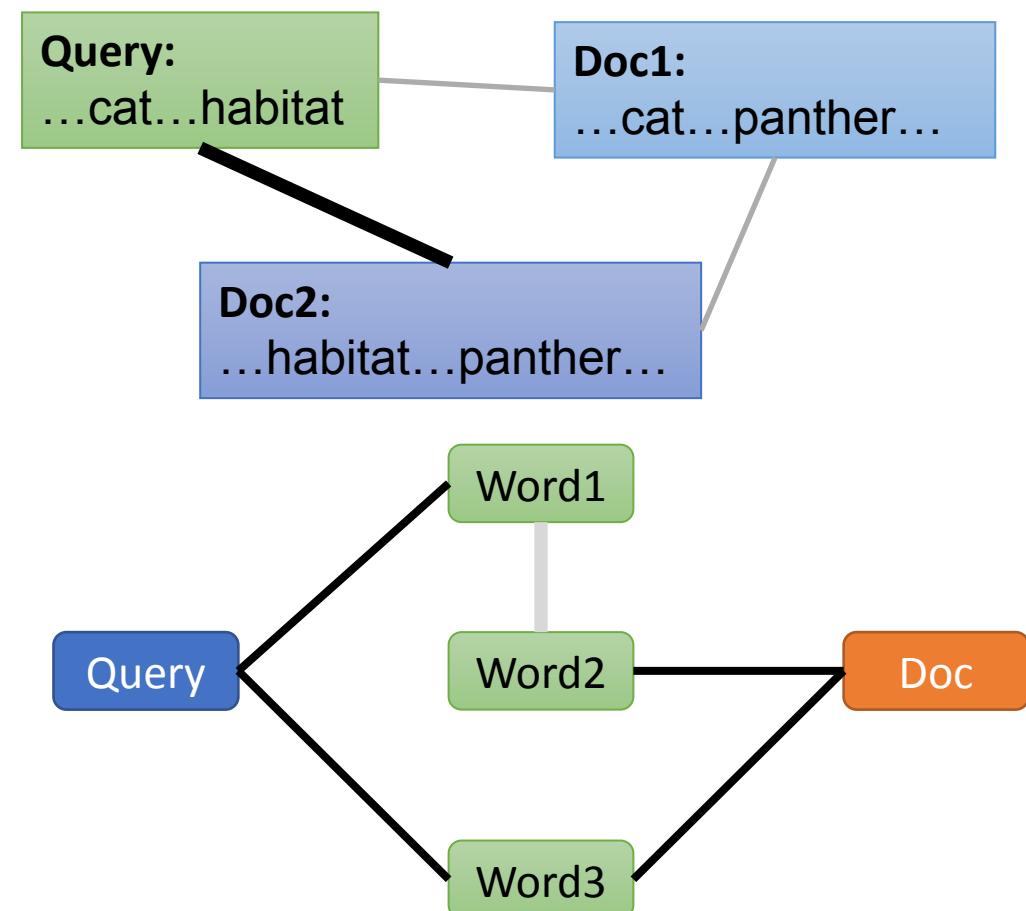
CS6101

Learning word and text representations

- Historical perspective
- Supervised
 - For ranking (this module)
- Unsupervised (BC)
 - Distributional vectors, word clustering, matrix factorization
- Unsupervised (AD)
 - Word2vec, GloVe (next)
- Contextualized
 - ConvNet, RNN, Transformer (later)
- Pretrained then fine-tuned

Pseudo-relevance feedback (PRF)

- Synonymy and polysemy
- Need to match documents to queries without any shared word
- Practical approach: [pseudo-relevance feedback \(PRF\)](#)
 - Process query from user to get top hits
 - Assume these are relevant
 - Extract keywords from these documents
 - Pad query (perhaps with smaller weight)
 - Process padded query
 - Return merged result lists
- Why stop at two queries?
- How to set magic weights?



Distributional vectors

- Finite labeled data, feature sparsity, and out-of-vocabulary (OOV) words/features have always troubled POS and NER tagging and estimating n-gram statistics
- E.g. we saw car in the training sequences but never sedan, or Shanghai in test data but only other cities in training data
- Unlabeled corpus has enough clues that these are related
- A well-established partial fix is word cluster features
- First [proposed by IBM researchers](#) in 1992
- Given a word like tezgüino, collect all context windows (say) at most 11 words wide centered on it
- From these contexts collect a bag of other words, count them
- Possibly transform from raw counts to TFIDF

- A bottle of tezgüino is on the table.
- Everybody likes tezgüino.
- Tezgüino makes you drunk.
- We make tezgüino out of corn.

What is tezgüino?

Word clusters

- Represent as a sparse vector in a space as large as the corpus vocabulary
 - perhaps scale to unit length
- This is the distributional vector for tezgüino, sort of a mini-doc
 - bottle, drunk, corn
- Turns out the d.v.'s of similar/related words are similar
 - tezgüino vs tequila vs toddy
- Can cluster these d.v.'s using standard clustering tools; see
https://en.wikipedia.org/wiki/Brown_clustering
 - Best number of clusters may be application-dependent
- In many NLP tasks, one of the features for each token is its cluster ID

Progression of text representation

- Sparse vectors, TFIDF, BM25
- Distributional vectors
- Dense representation of DVs (via PCA, say)
- Latent semantic indexing (SVD)
- Random walks on word-doc bipartite graph
- word2vec and GloVe
- Weighted average of words in passage
- 1d convolutional network
- RNNs, GRUs, LSTMs
- Transformers

PRF as associative retrieval

- If $r(i)$, $s(i, j)$ large, may also include j
 - Even if $r(j)$ is small
- Query {cat}, i includes {cat, panther}, j contains {panther} and is actually more relevant than i
- $r(j) < r(i)$ because of limitation of sparse matching (“lexical gap” between ‘cat’ and ‘panther’)
- How can we discover (ideally without supervision) that ‘cat’ and ‘panther’ are related?
 - Perpendicular axes in any sparse vector space model

A graph-based formulation

- Nodes i with ‘raw’ scores $r(i)$
- Undirected edges (i, j) with weights $s(i, j)$
- ‘Refined’ scores $y(i)$
- Two part objective
 - Want small $\sum_i (y(i) - r(i))^2$ (keep refined scores close to raw scores)
 - Want small $\sum_{(i,j)} s(i, j)(y(i) - y(j))^2$ (if i, j are similar, refined scores should be close)
- Convex optimization; solve for and sort by $y(i)$
- Above we assume $r(i)$ and $s(i, j)$ are fixed inputs

“Graph Laplacian”; we will return to this later

Training graph formulation

- For each query, we are given a subset of relevant nodes V_+ (remaining nodes V_\ominus)
 - Want $y(g) > y(b)$ for all $g \in V_+, b \in V_\ominus$
 - Usual hinge loss $[y(b) + m - y(g)]_+$
 - Three-part loss objective
 - If $r(i)$ and $s(i, j)$ are output from neural networks, gradient of loss can be backpropagated into these networks
 - Still does not explore word representation

Hardwired → trained ‘IDF’

- If we are given sufficient $\{\langle q, D_{q\oplus}, D_{q\ominus} \rangle\}$, can we learn something like IDF?
 - Or better? IDF is not informed by query word distribution
- Represent query q and doc x using normalized TF vectors alone, called $\mathbf{q}, \mathbf{x} \in \mathbb{R}^W$
- Instead of defining similarity as $\mathbf{q} \cdot \mathbf{x} = \mathbf{q}^\top \mathbf{x}$, define as $\mathbf{q}^\top \mathbf{A} \mathbf{x}$
- Where \mathbf{A} is a $W \times W$ non-negative diagonal matrix
 - E.g., define $\mathbf{A} = \exp(\mathbf{B})$ or $\mathbf{A} = \text{ReLU}(\mathbf{B})$ with unconstrained \mathbf{B}
- $A[w, w]$ represents the importance (similar to IDF) of word w
- Can try RankSVM-style constraint/loss

$$\mathbf{q}^\top \mathbf{A} \mathbf{x}_{q\oplus} + \xi_q \geq \mathbf{q}^\top \mathbf{A} \mathbf{x}_{q\ominus} + 1$$

Does this work?
Programming
assignment?

Structure of \mathbf{A}

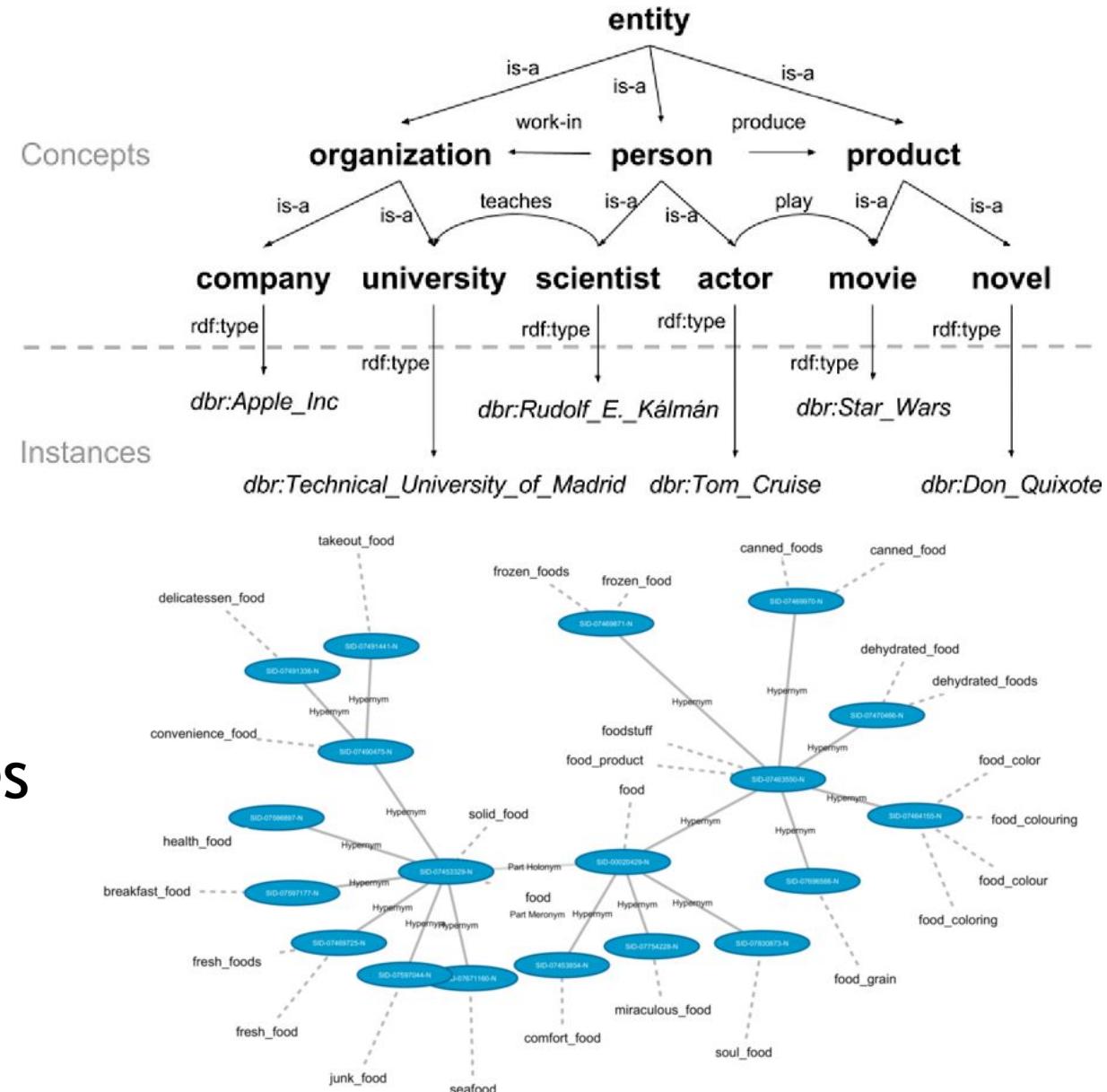
- Can we discover that ideally $A[\text{scold}, \text{rebuke}] > 0$?
 - Can improve recall
- We can relax \mathbf{A} to a general matrix
 - Should \mathbf{A} be symmetric? Diagonal-dominant? Sparse?
 - Do we have enough training data?
- Which $A[i, j]$ do we expect to be large?

Constraining A

- One way is to model

$$A = \mathbb{I} + \lambda \text{ReLU}(B^\top B)$$

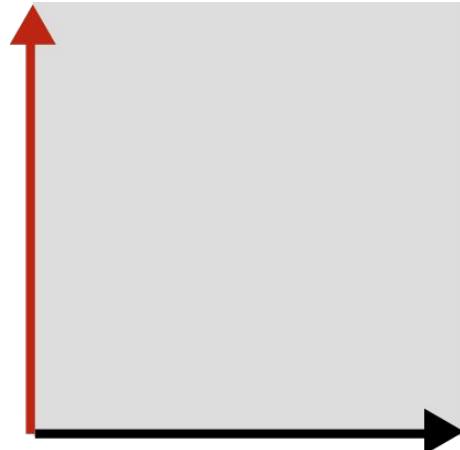
- where $B \in \mathbb{R}^{M \times W}$ (low-rank)
- and λ is kept small
- Another way is to use lexical networks like wordnet
 - Allow $A[i, j] > 0$ only if words i, j are connected in wordnet by some restricted types of edges (synonym, hypo/hypernym)



Finishing up

- In any case, similarity $s_B(q, x_i)$
- Now the loss is pairwise
 - Aka “triplet loss”, the third element being q
- $\sum_q \sum_{g,b} \text{ReLU}(\text{margin} + s_B(q, b) - s_B(q, g))$
 - May need to tune margin here
- Minimize pairwise loss wrt B (or sparse A) to learn translation matrix (with possibly “IDF” on diagonals)
- Stare a bit at $B \in \mathbb{R}^{M \times W}$
 - Looks like M -dim embedding of each word
 - Fitted for end-task of retrieval

Linear transformations and SVD



A square with orienting arrows

- A: stretched on black axis
- B: compressed on black axis
- C: rotated
- D: reflected
- E: sheared

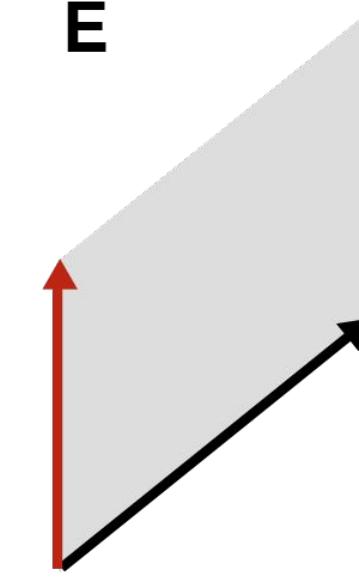
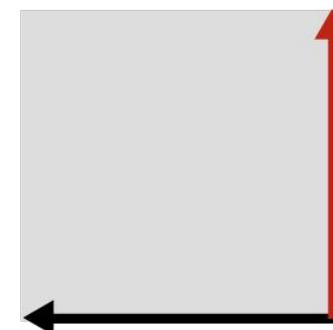
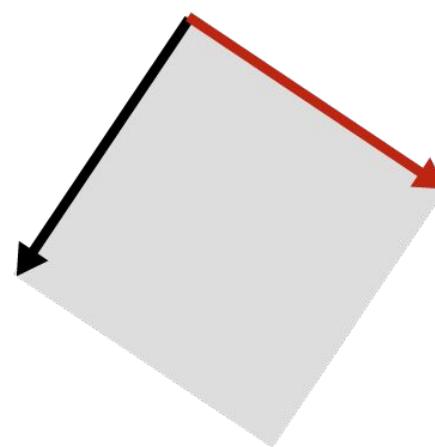
A

B

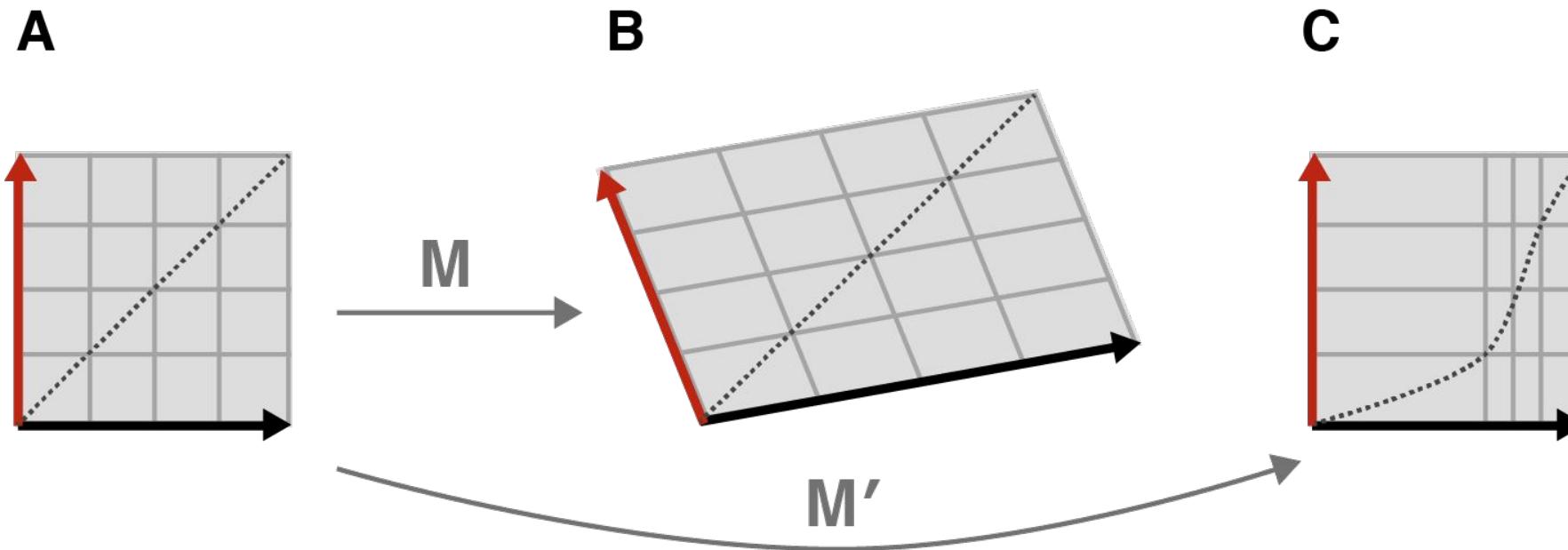
C

D

E



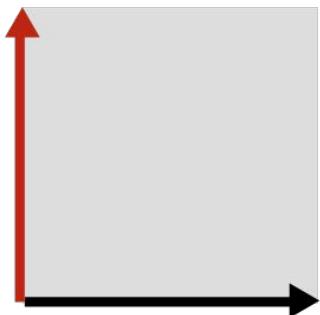
Let M be a linear transformation



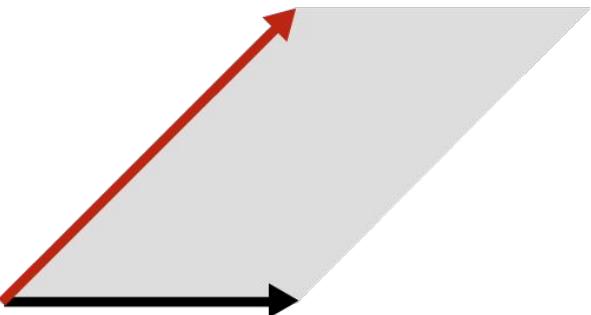
- A: original square
- B: under linear transformation M
 - Linear transformation transforms a straight line to a straight line
- C: under non-linear transformation M'

Rotate, then stretch/compress/reflect

A

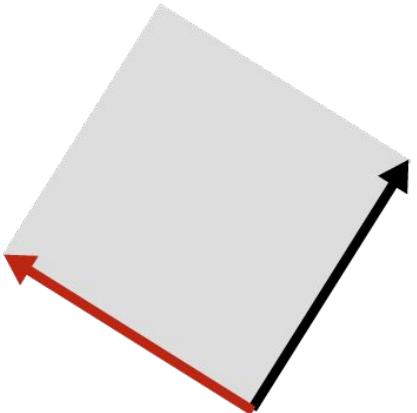


M

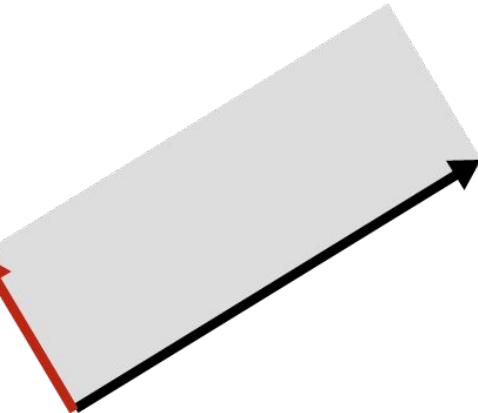


- Under M , original square is sheared
- But rotated square is only stretched or compressed

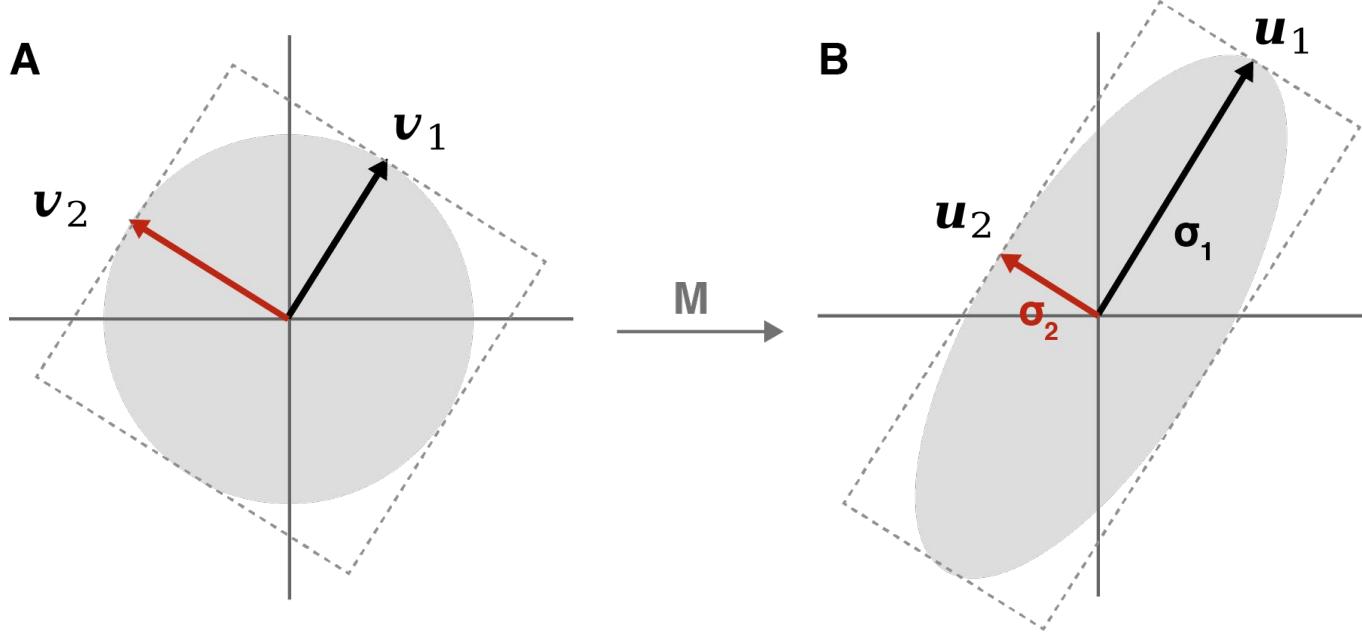
B



M



Circle to ellipse



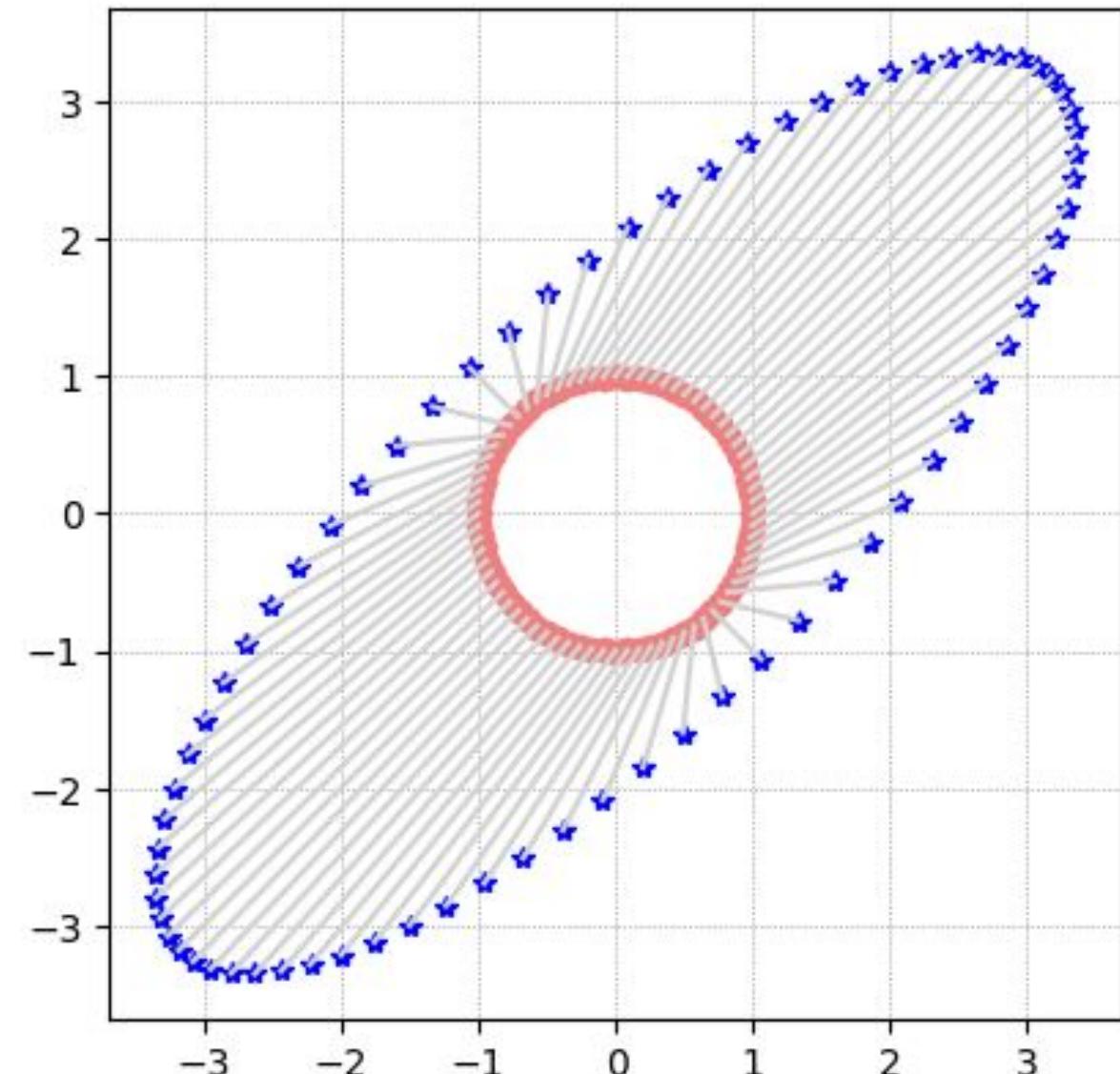
$$Mv_1 = \sigma_1 u_1$$
$$Mv_2 = \sigma_2 u_2$$

Where u_1, u_2
are unit vectors

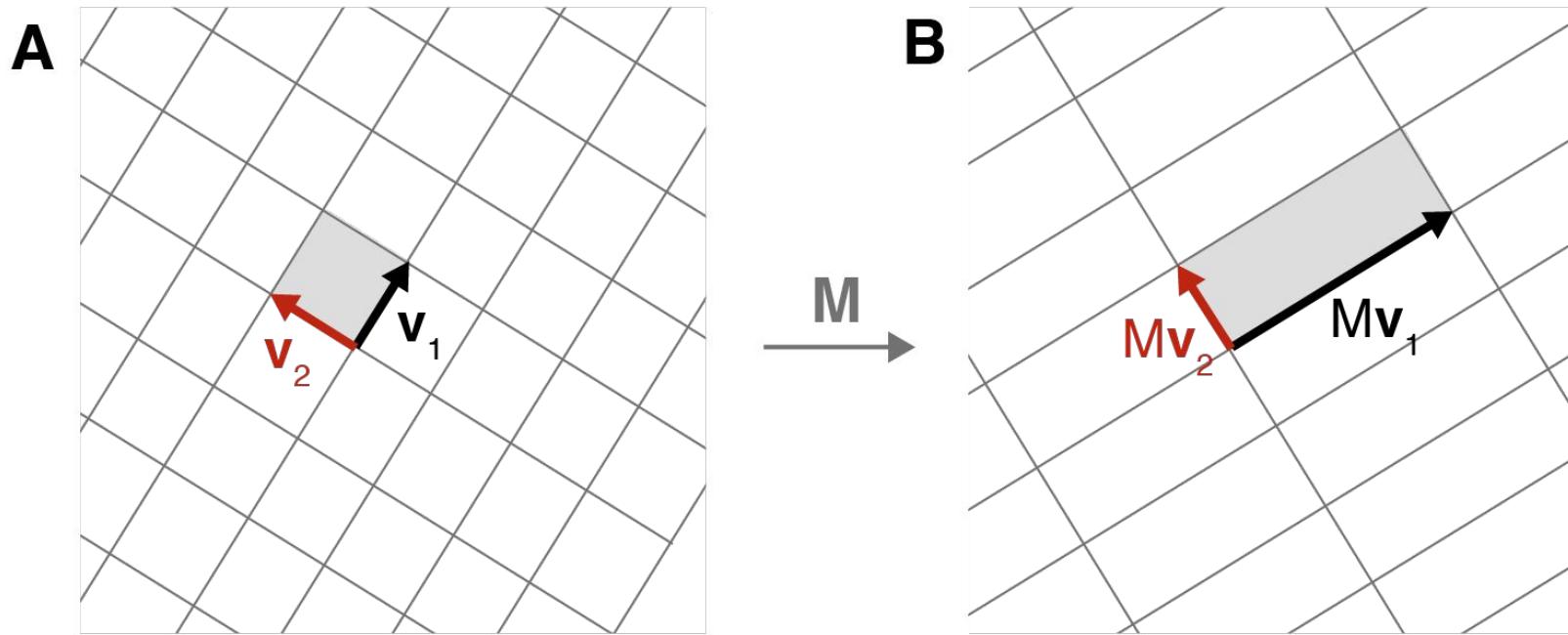
- Inscribe a circle inside the rotated square
- Under M , circle maps to ellipse
- Major and minor axes define stretch/shrink
 - Black and red directions do not change, only length

Effect of matrix transform on unit circle

- $M = \begin{bmatrix} 5 & 1.5 \\ 1.5 & 3 \end{bmatrix}$
- Circle mapped to ellipse
- Two radial lines (major and minor axes) get scaled but not turned
 - Eigenvectors of M
- Also explains why power iteration finds largest eigenvector



In other words



- Domain spanned by (unit) basis vectors $\mathbf{v}_1, \mathbf{v}_2$
- Pre-rotate it through a specific angle, then apply linear transformation M
- Amounts to transformation by a diagonal matrix

Transforming any vector \mathbf{x}

- Given (unit) basis vectors $\mathbf{v}_1, \mathbf{v}_2$

- Any vector \mathbf{x} can be written as

$$\mathbf{x} = (\mathbf{x} \cdot \mathbf{v}_1)\mathbf{v}_1 + (\mathbf{x} \cdot \mathbf{v}_2)\mathbf{v}_2$$

- After transformation by M , we get

$$M\mathbf{x} = (\mathbf{x} \cdot \mathbf{v}_1)M\mathbf{v}_1 + (\mathbf{x} \cdot \mathbf{v}_2)M\mathbf{v}_2$$

- Which can be rewritten as

$$M\mathbf{x} = (\mathbf{x} \cdot \mathbf{v}_1)\sigma_1\mathbf{u}_1 + (\mathbf{x} \cdot \mathbf{v}_2)\sigma_2\mathbf{u}_2$$

- Rearrange to

$$M\mathbf{x} = \mathbf{u}_1\sigma_1(\mathbf{x} \cdot \mathbf{v}_1) + \mathbf{u}_2\sigma_2(\mathbf{x} \cdot \mathbf{v}_2)$$

- Since this is true for all \mathbf{x} , we get

$$M = \mathbf{u}_1\sigma_1\mathbf{v}_1^\top + \mathbf{u}_2\sigma_2\mathbf{v}_2^\top$$

$$M\mathbf{v}_1 = \sigma_1\mathbf{u}_1$$
$$M\mathbf{v}_2 = \sigma_2\mathbf{u}_2$$

Generalizing to more dimensions

$$M = U \Sigma V^*$$

The diagram shows the decomposition of a matrix M into U , Σ , and V^* . The matrix M is represented by a large gray rectangle. An equals sign follows it. To the right of the equals sign, there is a vertical stack of three rectangles. The first rectangle, labeled U below it, has two columns: a solid column on the left and a dashed column on the right. A diagonal band of gray shading runs from the top-left of the solid column to the bottom-right of the dashed column. The second rectangle, labeled Σ below it, is a square with a diagonal band of gray shading running from the top-left to the bottom-right. The third rectangle, labeled V^* below it, is a single column with a solid gray rectangle at the top and a white rectangle at the bottom.

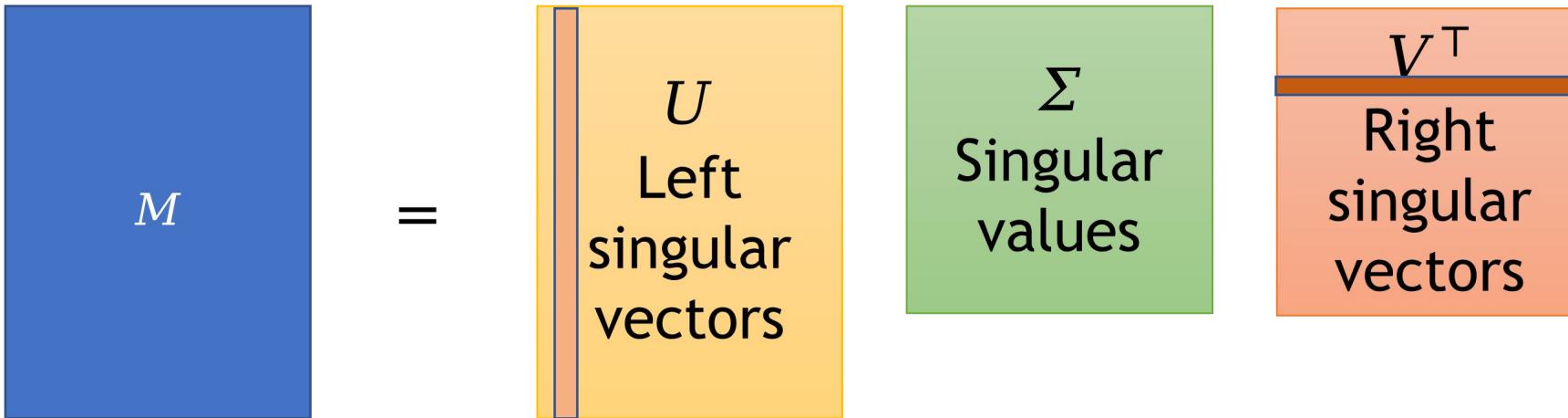
- In 2d, we got $M = [\mathbf{u}_1 \ \mathbf{u}_2] \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix} [\mathbf{v}_1^\top \ \mathbf{v}_2^\top]$
- Written compactly as $M = U \Sigma V^\top$
- $\Sigma = \text{diag}(\{\sigma_i\})$ with $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$
- Columns of U are unit vectors, perpendicular to each other
- So are columns of V

Signal to noise

- Rank-revealing decomposition

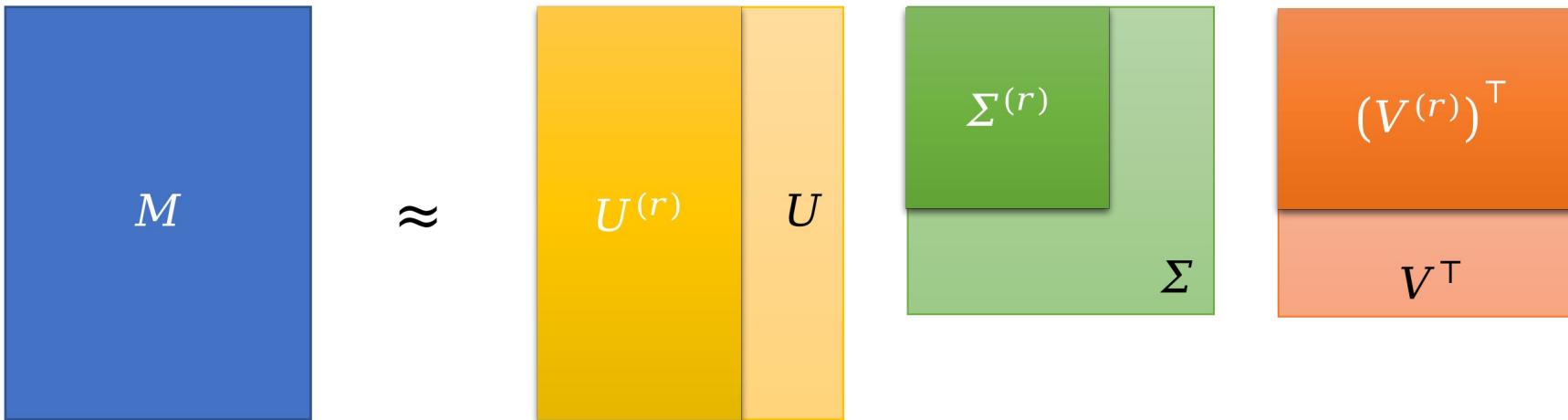
- If $\text{rank}(M) = R$, then $\sigma_{R+1} = \sigma_{R+2} = \dots = 0$
- Frobenius norm of a matrix $\|M\|_F^2 = \sum_{i,j} M_{i,j}^2$
 - “Total energy” in a matrix
- Trace of a square matrix $\text{tr}(A) = \sum_i A_{i,i}$ is the sum of diagonal elements
- $\|M\|_F^2 = \text{tr}(MM^\top) = \text{tr}(U\Sigma V^\top V\Sigma^\top U^\top) = \text{tr}(U\Sigma\Sigma^\top U^\top) = \text{tr}(\Sigma\Sigma^\top U^\top U) = \text{tr}(\Sigma\Sigma^\top) = \sum_i \sigma_i^2$
- “Total energy” of matrix = sum of squares of singular values

Low rank approximation



- Which matrix B with rank $r \leq R$ is the “closest approximation” to input matrix M ?
- I.e., minimize $\|M - B\|_F^2$ s.t. $\text{rank}(B) = r$

Low rank approximation



- Which matrix B with $\text{rank } r \leq R$ is the “closest approximation” to input matrix M ?
- I.e., minimize $\|M - B\|_F^2$ s.t. $\text{rank}(B) = r$
- Solution is $B = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$
- With error $\|M - B\|_F^2 = \sum_{i=r+1}^R \sigma_i^2$ = the “lost power”
- What r is sufficient for doc-term matrices? What is approx rank?

Connection to eigen system

- $M = U \Sigma V^\top$
- $C = M M^\top = U \Sigma V^\top V \Sigma^\top U^\top = U \Sigma \mathbb{I} \Sigma^\top U^\top = U \Sigma^2 U^\top$
- $CU = U \Sigma^2 U^\top U = U \Sigma^2 \mathbb{I} = \Sigma^2 U$
- $C\mathbf{u}_1 = \sigma_1^2 \mathbf{u}_1, C\mathbf{u}_2 = \sigma_2^2 \mathbf{u}_2$
 - Left singular vectors of M are eigenvectors of $M M^\top$
 - Squares of singular values are eigenvalues
- Similar properties for right singular vectors and $M^\top M$

Mutually recursive similarity

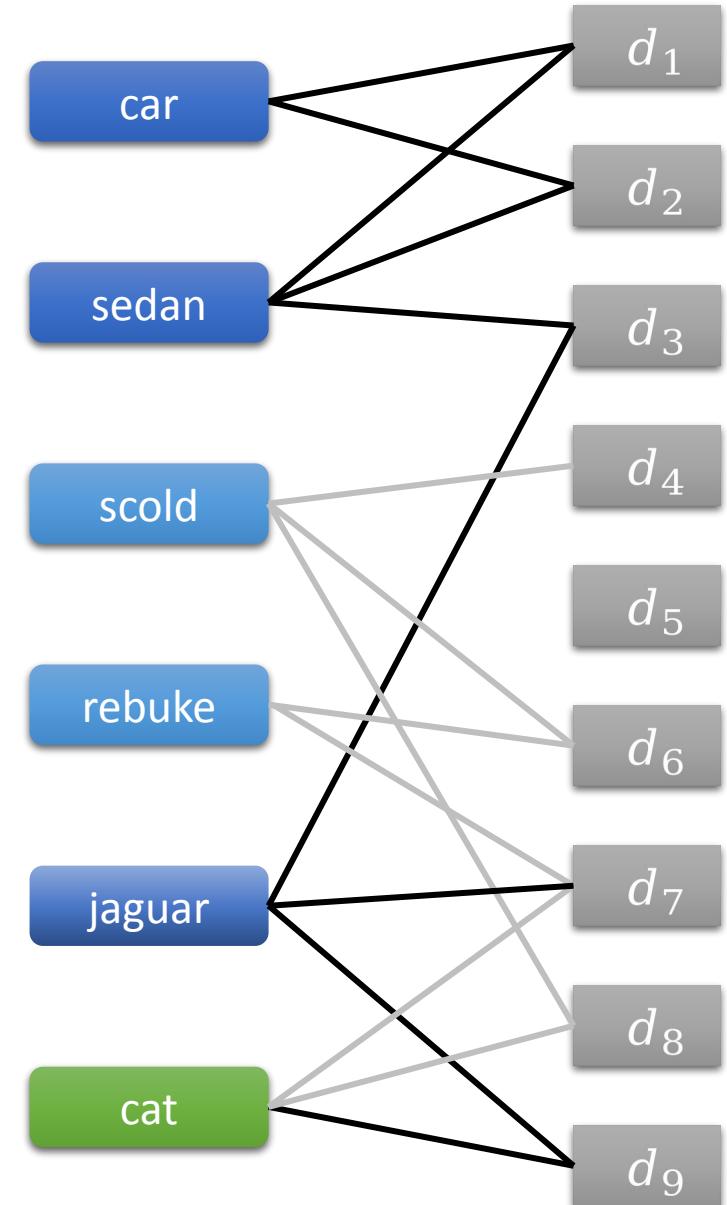
- Two words are similar if they appear in similar documents
- Two documents are similar if they contain similar words
- Turn this into a random walk problem
- Bipartite graph between words and docs
- Expressed via adjacency matrix $E \in \{0,1\}^{W \times D}$

Back to term-doc matrix

- Suppose M is the term (rows) \times doc (cols) matrix
- What is $(MM^\top)[i, j]$?
 - Number of docs that contain both words i and j
 - The number of 2-hop paths between bipartite term-document graph – a form of similarity
- Closely related to principal component analysis (PCA)

Random walk on word-doc graph

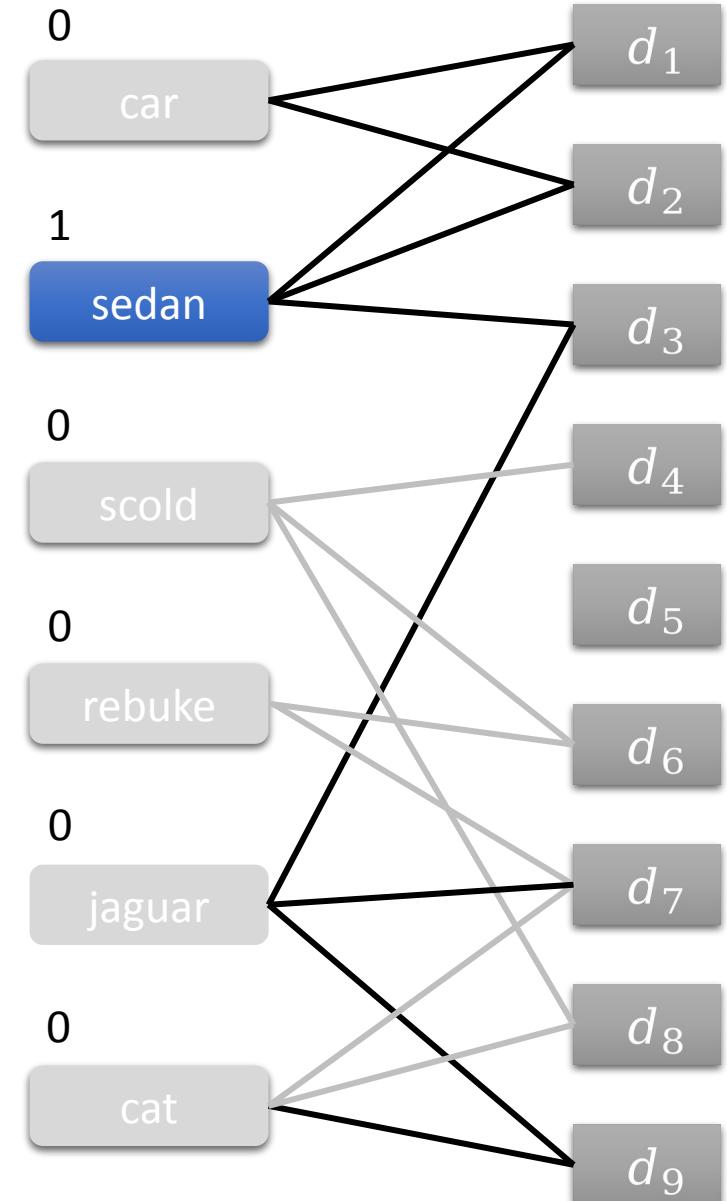
- Bipartite graph between words and docs
- Expressed via adjacency matrix $E \in \{0,1\}^{W \times D}$
- Two similar questions
 - If a random surfer starts from w_i , and resets to w_i w.p. α every step, what is the steady state probability of visiting w_j ?
 - How many paths of length 2, 4, 6,... exist between words w_i, w_j ?
- $EE^T \in \mathbb{Z}_+^{W \times W}$; $(EE^T)[i, j]$ gives the number of 2-hop paths between w_i, w_j
- $(EE^T)^2[i, j]$ gives the number of 4-hop paths between w_i, w_j , and so on



Spreading activation

- $(EE^\top)^k[i, j]$ gives the number of $2k$ -hop paths between w_i, w_j
- Start with a vector $\mathbf{u}^{(0)} \in \mathbb{R}^W$
 - Elements sampled from some distribution
 - Or the impulse at word w_0 : set $\mathbf{u}^{(0)} = \delta_{w_0}$
- Repeat
 - $\mathbf{v}^{(k+1)} \leftarrow (EE^\top)\mathbf{u}^{(t)}$
 - $\mathbf{u}^{(k+1)} \leftarrow \mathbf{v}^{(k+1)} / \|\mathbf{v}^{(k+1)}\|_1$
- (If this converges) finds an eigenvector \mathbf{u} such that $(EE^\top)\mathbf{u} = \lambda \mathbf{u}$

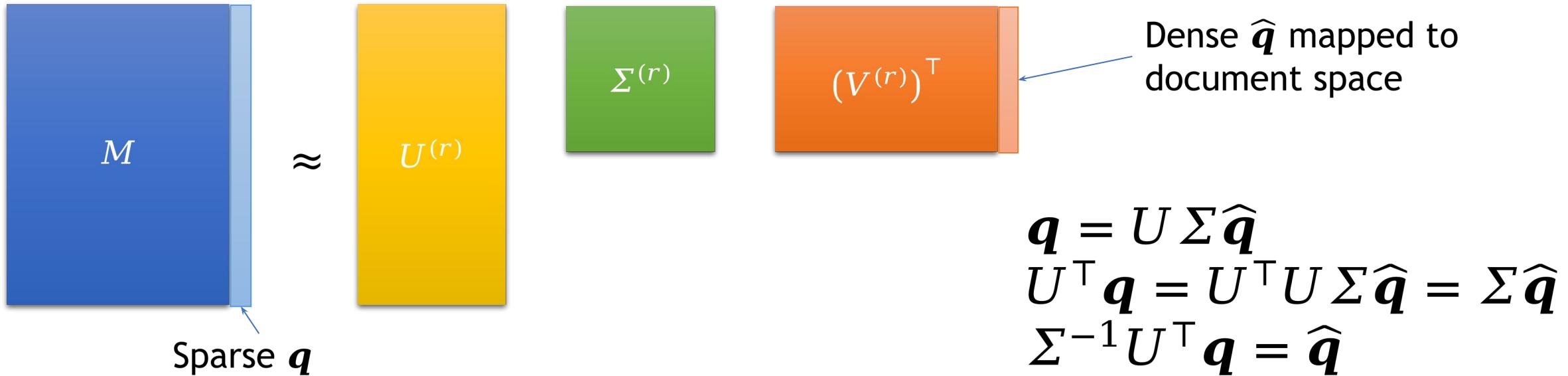
Depending on whether E is connected, $\mathbf{u}^{(\infty)}$ may or may not be sensitive to $\mathbf{u}^{(0)}$. More about this later.



Using SVD in search

- Sparse term-doc matrix M decomposed as $U \Sigma V^T$
- Usually truncated to ≈ 300 dimensions
- Each row of U thus gives a 300-dim representation of a word
- Each row of V gives a 300-dim representation of a doc
- These “embeddings” are no longer sparse
- “Latent semantic indexing” or LSI, 1988
- Query is like a new sparse doc vector

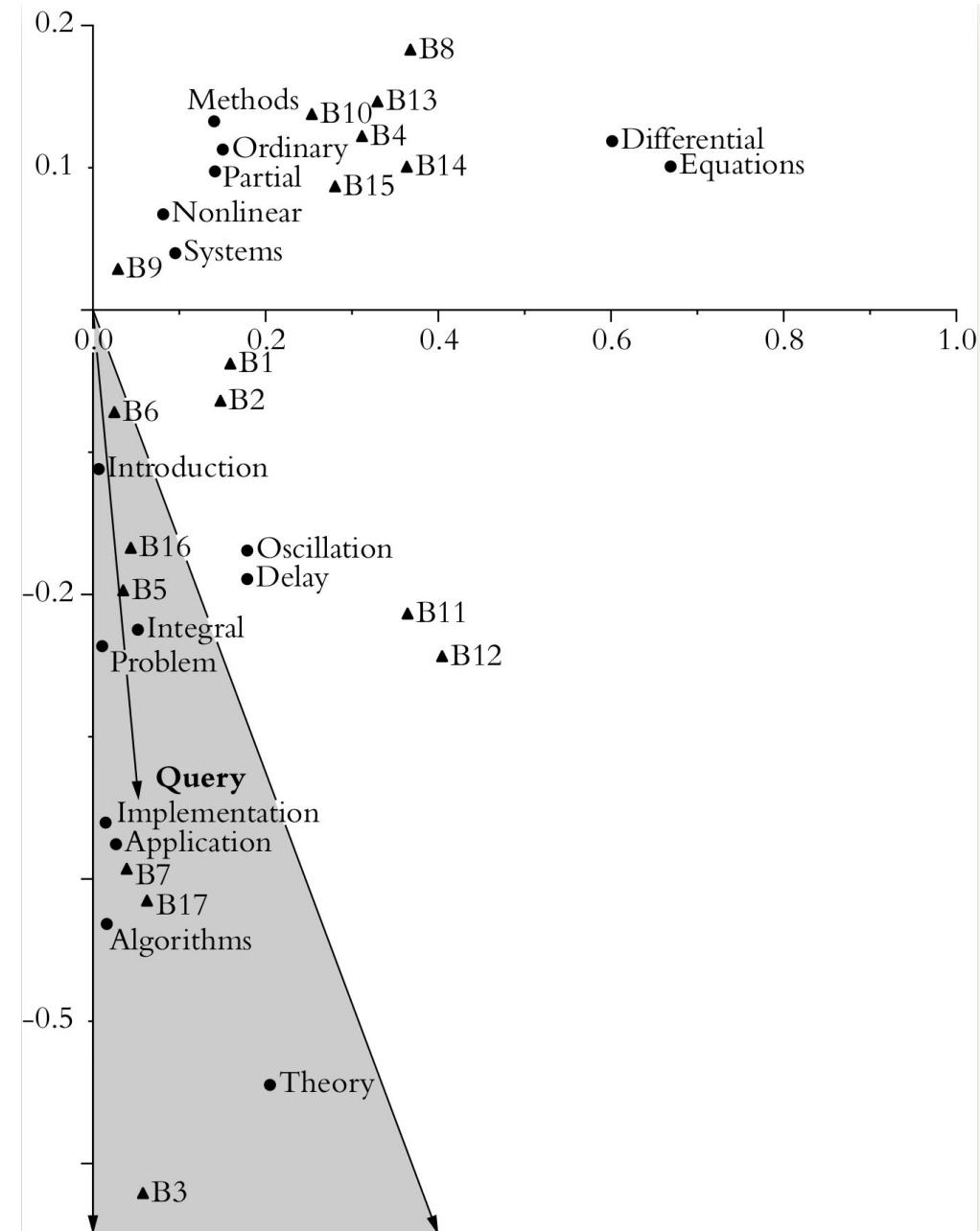
LSI in search, continued



- Pad M with sparse query vector \mathbf{q}
 - Assume perturbation small enough that U, Σ, V do not change much
 - Project sparse query \mathbf{q} to dense $\widehat{\mathbf{q}}$ in doc space
 - Find nearby doc vectors as $\max_K \widehat{\mathbf{q}} \cdot \mathbf{v}_j$ (cannot use inverted index)

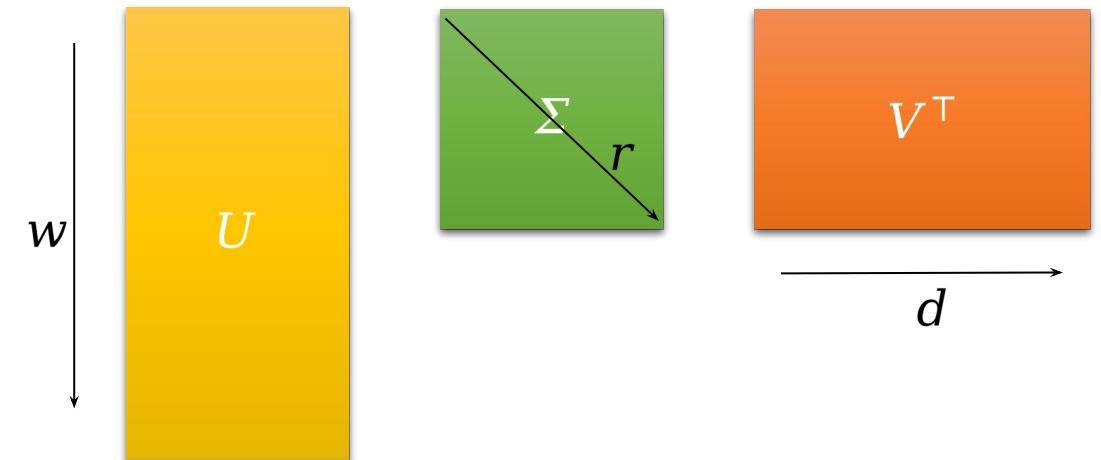
LSI example

Label	Titles
B1	A Course on <u>Integral Equations</u>
B2	Attractors for Semigroups and Evolution Equations
B3	Automatic Differentiation of <u>Algorithms</u> : <u>Theory</u> , <u>Implementation</u> , and <u>Application</u>
B4	Geometrical Aspects of <u>Partial Differential Equations</u>
B5	Ideals, Varieties, and <u>Algorithms</u> –An <u>Introduction</u> to Computational Algebraic Geometry and Commutative Algebra
B6	<u>Introduction</u> to Hamiltonian Dynamical <u>Systems</u> and the <u>N-Body Problem</u>
B7	Knapsack Problems: <u>Algorithms</u> and Computer <u>Implementations</u>
B8	Methods of Solving Singular <u>Systems</u> of <u>Ordinary Differential Equations</u>
B9	<u>Nonlinear Systems</u>
B10	<u>Ordinary Differential Equations</u>
B11	<u>Oscillation Theory</u> for Neutral <u>Differential Equations</u> with <u>Delay</u>
B12	<u>Oscillation Theory</u> of Delay <u>Differential Equations</u>
B13	Pseudodifferential Operators and <u>Nonlinear Partial Differential Equations</u>
B14	Sync <u>Methods</u> for Quadrature and <u>Differential Equations</u>
B15	Stability of Stochastic <u>Differential Equations</u> with Respect to Semi-Martingales
B16	The Boundary <u>Integral Approach</u> to Static and Dynamic Contact Problems
B17	The Double Mellin-Barnes Type <u>Integrals</u> and Their <u>Applications</u> to Convolution <u>Theory</u>



Non-negative matrix factorization

- Unconstrained SVD $M = U \Sigma V^\top$
- Rank-constrained approx $M \approx U^{(r)} \Sigma^{(r)} (V^{(r)})^\top$
- A different way to constrain:
 - All elements of U, Σ, V must be non-negative
 - Each column of U and V sum to 1
 - Diagonal elements of Σ sum to 1
- Multinomial distribution over the joint (w, d) space
 - Both regarded as random variables



- For each r , $\sum_w U[w, r] = 1$
- Think $\sum_w p(w|r) = 1$
- For each r , $\sum_d V[d, r] = 1$
- Think $\sum_w p(d|r) = 1$
- $\sum_r \Sigma[r] = 1 \rightarrow \sum_r p(r) = 1$

Dyadic or aspect model or PLSI

- Joint distribution $p(w, d)$ not product of marginals
 - I.e., w, d are strongly dependent
- Bottleneck through topics r
- $p(w, d) = \sum_r p(r)p(w, d|r)$ is always true
- We approximate $p(w, d|r) \approx p(w|r)p(d|r)$
 - “Word w and doc d are conditionally independent given topic r ”
- Overall $p(w, d) \approx \sum_r p(r)p(w|r)p(d|r)$
 - (Looks structurally similar to SVD)

Fitting a dyadic model

- Model parameterized as $p(r) = \Sigma[r]$, $p(w|r) = U[w, r]$, $p(d|r) = V[d, r]$
- Total $WR + R + DR$ real parameters
- Event probability $p(w, d) \approx \sum_r p(r) p(w|r)p(d|r)$
- Let the term-doc matrix contain (w, d) counts $n[w, d]$
- Objective is $\max_{U, \Sigma, V} \sum_{w,d} n[w, d] \log p(w, d)$
- Subject to U, Σ, V satisfying probability-like constraints

EM optimization

- Given current (“M variables”) $p(r), p(w|r), p(d|r)$, alternate E and M steps
 - Estimate posterior (E)

$$p(r|w, d) = \frac{p(r)p(w, d|r)}{\sum_{r'} p(r')p(w, d|r')} = \frac{p(r)p(w|r)p(d|r)}{\sum_{r'} p(r')p(w|r')p(d|r')}$$

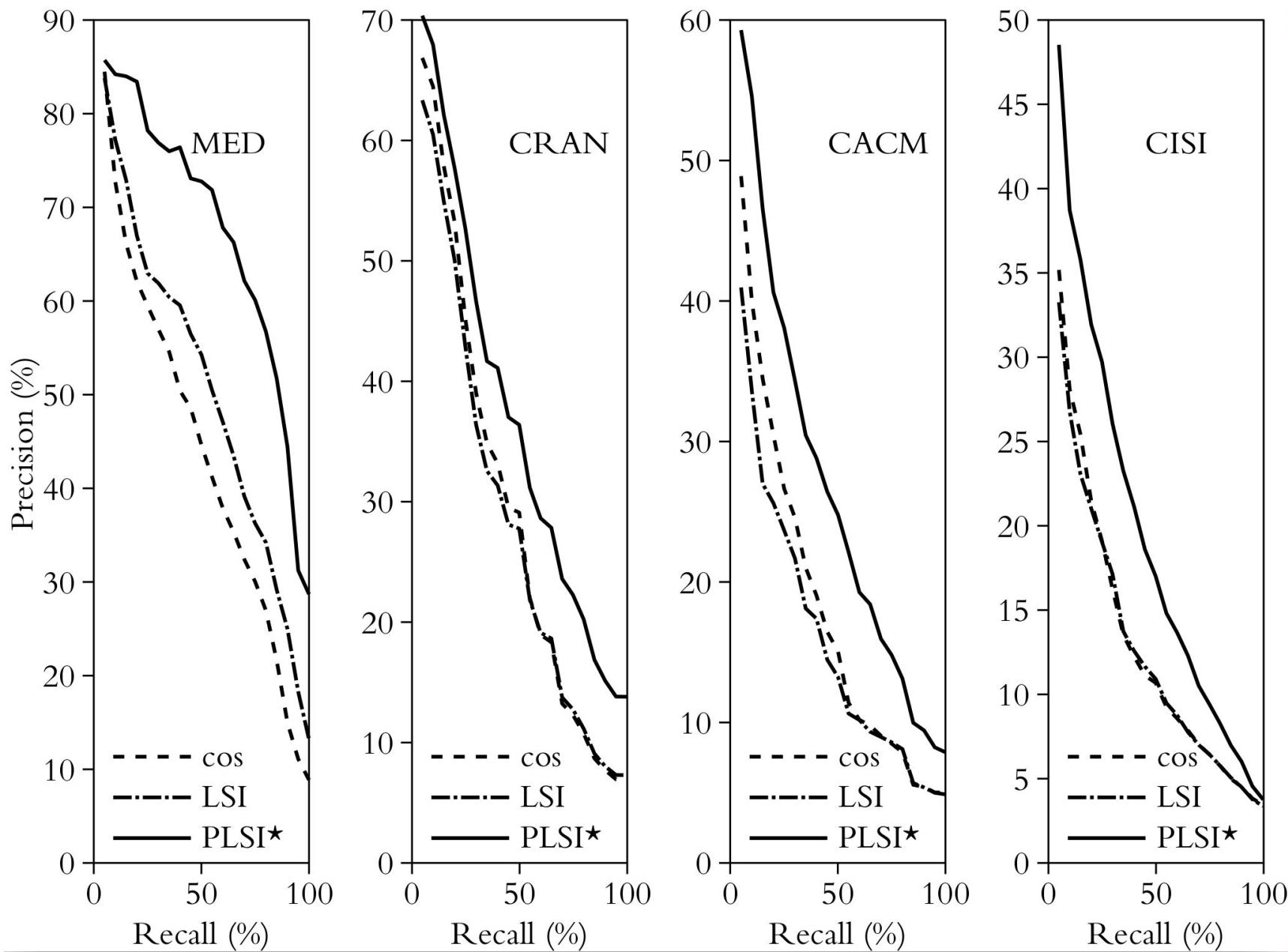
- Compute next “M variable” estimates

$$p(r) = \frac{\sum_{w,d} n(w, d) p(r|w, d)}{\sum_{r'} \sum_{w,d} n(w, d) p(r'|w, d)}$$

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

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

TFIDF vs LSI vs PLSI



Summary

- Lexical gap has been known for ages
- Pre-deep learning approaches
 - Pseudo relevance feedback
 - Word clustering
- Approaching word representation
 - Distributional vectors
 - LSI, PLSI
- Next: neural contextualized word rep.