

DMML 9 Apr 2019

Informational Retrieval & Web Search

Page Rank

Given a query/topic

Retrieve a collection of relevant documents

Classify these as authoritative sources of information vs sources of links to information

Jon Kleinberg (same year as Page Rank)

Hubs

Authorities

Stars & Directors

Likewise hubs get value from the  
quality of authorities they point to

& vice versa

document  $i$

$$h(i) = \sum_j a(j) W[i, j]$$

hub score

$$W_{\text{Docs}} \begin{bmatrix} \text{Docs} \\ \text{---} | \end{bmatrix}$$

Symmetrically

$$a(i) = \sum_j h(j) W[j, i]$$

$$H = WA$$

$$A = W^T H$$

$$H = \underline{\underline{WW^T}} H$$

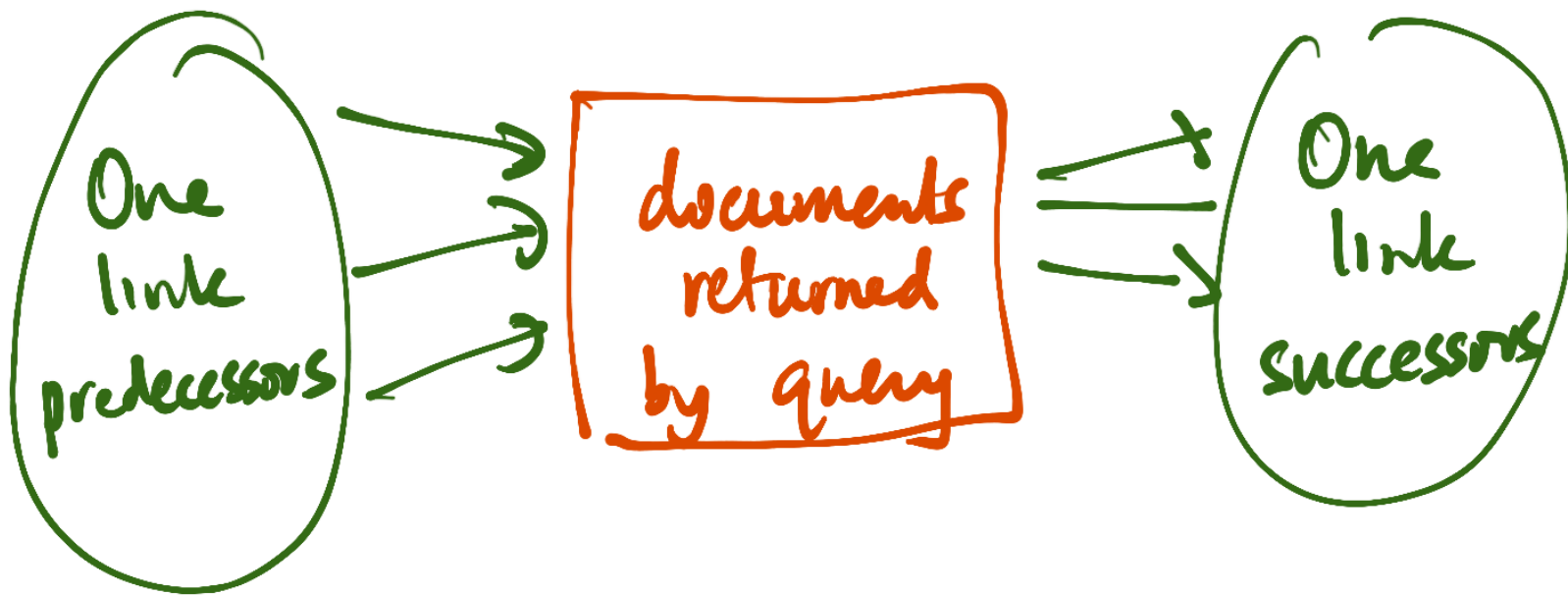
$$A = \underline{\underline{W^T W}} A$$

Recall PageRank

$$\pi = M\pi$$

What is  $W$ ? Which documents?

query  $\rightarrow$  relevant documents



Space of  $W$  — control size

Control size

200 documents in starting set

Limit how many new documents each  
one "pulls in" (50)

$W$  is the incidence matrix of this set

$L$  need not be irreducible, aperiodic

We may have multiple "largest" eigenvalues

L indicative of a partition (not irreducible)

Ambiguous words

Jaguar

Two "communities"



Kleinberg called this HITS

Hypertext Induced Topic Search

---

Changing gears

Two problems with search

Words have different meanings

Multiple words have same meaning



# Vector Space

Term Doc Matrix with weights

Queries & documents as vectors and  
compare using cosine similarity

# Singular Value Decomposition (SVD)

$D$  is our term-doc matrix

$$\begin{array}{c} \text{terms} \\ \downarrow \\ m \end{array} \left[ \begin{array}{c} n \leftarrow \text{docs} \end{array} \right]$$

Unique way to write

$$D = U \Sigma V^T \quad \text{s.t.}$$

$m \times n \quad m \times m \quad m \times n \quad n \times n$

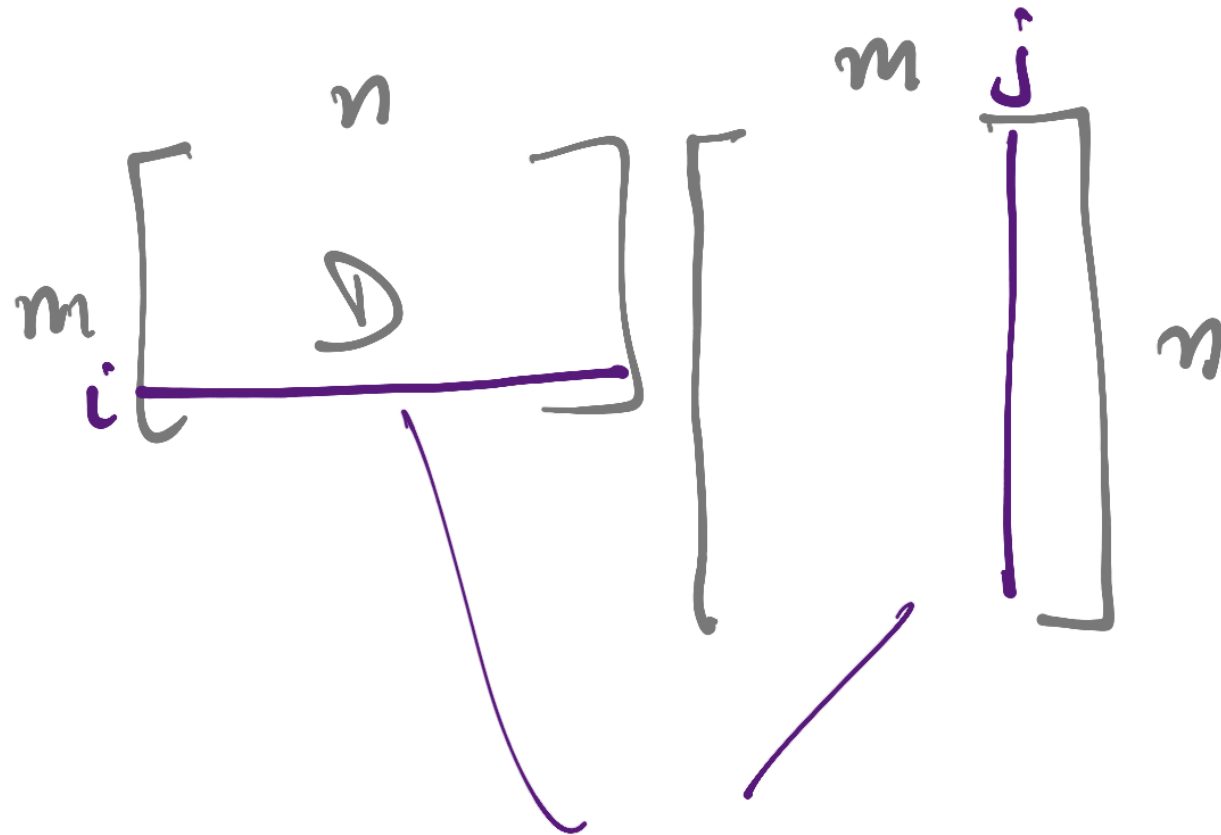
$U$  eigenvectors of  $DD^T$ ,  $U^T U = I$   
 $V$  eigenvectors of  $D^T D$ ,  $V^T V = I$

$\Sigma$  diagonal matrix  
 $\sigma_i = \sqrt{\lambda_i}$ ,  
 $m \begin{bmatrix} \sigma_1 & \sigma_2 & \dots & 0 & 0 \\ 0 & \sigma_1 & \dots & \sigma_r & 0 \end{bmatrix}$

eigenvalues of  $D$   
 in decreasing order

"singular values"

$U$  eigenvalues of  $DD^T$



Sum up all cooccurrences of  
 $t_i, t_j$

$$V = D^T D$$

- similarity across documents

$D$

$(\mathbf{d}_j)$

$\downarrow$

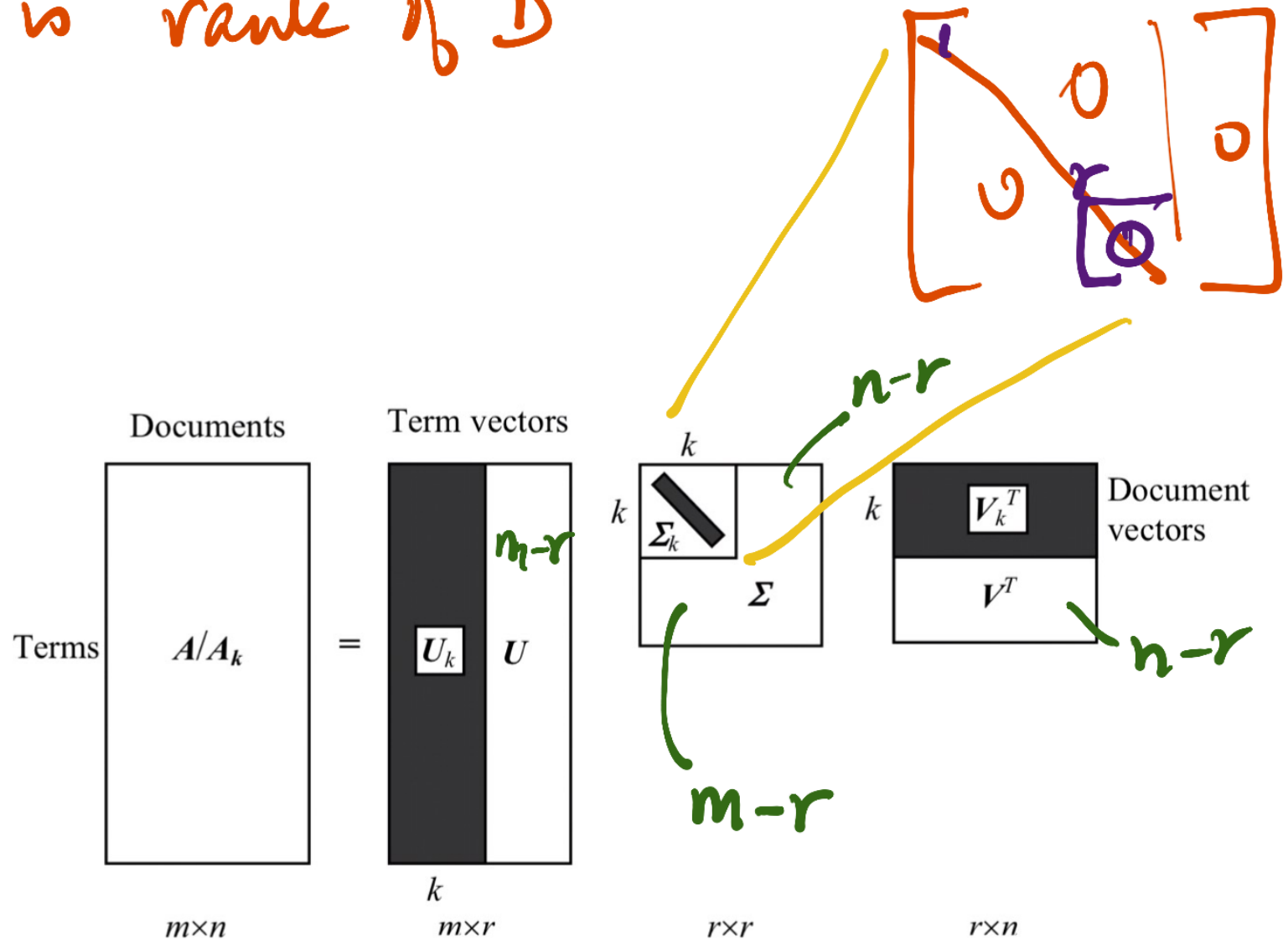
$(\mathbf{t}_i^T) \rightarrow$

$$\begin{bmatrix} x_{1,1} & \dots & x_{1,j} & \dots & x_{1,n} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{i,1} & \dots & x_{i,j} & \dots & x_{i,n} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{m,1} & \dots & x_{m,j} & \dots & x_{m,n} \end{bmatrix}$$



$\Sigma$  has  $r$  singular values - eigenvalues of  $D$

$r$  is rank of  $D$



Simplify

$$U \text{ to } m \times r$$

$$\Sigma \text{ to } r \times r$$

$$V \text{ to } r \times n$$



$$\begin{array}{ccccc}
 & U & & \Sigma & & V^T \\
 & & & & & (\hat{\mathbf{d}}_j) \\
 & & & & & \downarrow \\
 (\hat{\mathbf{t}}_i^T) \rightarrow & \left[ \begin{bmatrix} \mathbf{u}_1 & \dots & \mathbf{u}_l \end{bmatrix} \right] & \cdot & \begin{bmatrix} \sigma_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_l \end{bmatrix} & \cdot & \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_l \end{bmatrix}
 \end{array}$$

$\hat{t}_i$  is transformed version of  $t_i$

$\hat{d}_j$  is transformed version of  $d_j$

Symmetric matrix  $M$

Eigenvectors are orthogonal

$$\begin{array}{ccc} \lambda_1 & \lambda_2 & \lambda_3 \\ v_1 & v_2 & v_3 \end{array}$$

$$M \in 3 \times 3$$

Any vector  $u = (u_1, u_2, u_3)$

can be expressed as

$$n_1 v_1 + n_2 v_2 + n_3 v_3$$

$$M \cdot u$$

$$= M (n_1 v_1 + n_2 v_2 + n_3 v_3)$$

Eigenvector  $M v_i = \lambda_i v_i$

$$= (\lambda_1 v_1) n_1 + (\lambda_2 v_2) n_2 + \underline{\underline{(\lambda_3 v_3) n_3}}$$

Suppose  $\lambda_1 > \lambda_2 > \underline{\underline{\lambda_3}}$   
small

Ignore  $\lambda_3 v_3$ ,  
approximation of  $M \cdot u$

# Approximation of $M$

$D$  - restrict to  $k$  most significant  
terms (directions)

document

rank  $k$

keep only top  $k$  eigenvalues

Best approx is to truncate  $\Sigma$  to  $\Sigma_k$

$$D = U \Sigma V^T$$

$m \times r$     $r \times r$     $r \times n$



$$D' =$$

$m \times k$     $k \times k$     $k \times n$

$$\sum (D_{ij} - D'_{ij})^2$$

k-approx of SVD

How?

$$\begin{array}{ccccc} & U & & \Sigma & & V^T \\ & & & & & (\hat{\mathbf{d}}_j) \\ & & & & & \downarrow \\ (\hat{\mathbf{t}}_i^T) \rightarrow & \left[ \begin{bmatrix} \mathbf{u}_1 \end{bmatrix} \dots \begin{bmatrix} \mathbf{u}_l \end{bmatrix} \right] & \cdot & \begin{bmatrix} \sigma_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_l \end{bmatrix} & \cdot & \left[ \begin{bmatrix} \mathbf{v}_1 \end{bmatrix} \\ \vdots \\ \begin{bmatrix} \mathbf{v}_l \end{bmatrix} \end{bmatrix} \end{array}$$

What is  $\hat{\mathbf{d}}_j$ ?

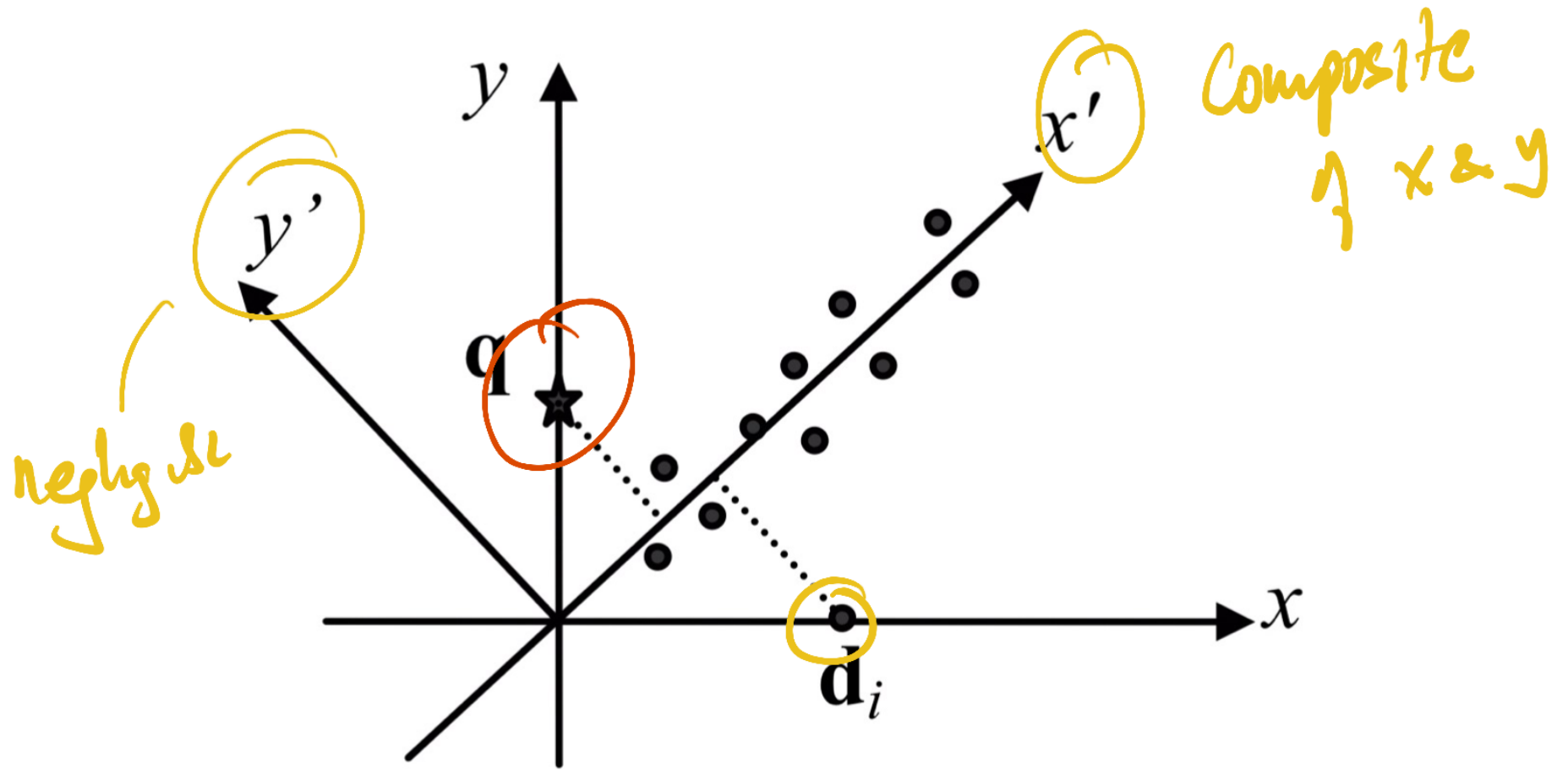
$$D = U_k \Sigma_k V_k^T \quad \leftarrow \text{after } k \text{ approx.}$$

Want  $V_k^T$  - representation of documents

$$U^T D = \underbrace{U^T U}_I \Sigma V^T$$

$$\Sigma^{-1} U^T D = V^T$$

$\downarrow$   
 $\frac{1}{\sigma_i}$



Rotate a query similarly

$$\tilde{q} = \Sigma^{-1} U^T q$$



Compare  $\hat{q}$  &  $\hat{d}_j$  as vectors

Latent Semantic Indexing