

Vector space reduction

As already said in the text-based ranking's notes, computing the similarity between two docs (cosine) in vector space is costly due to the high-dimensionality, therefore we would like to "pack" our vectors in fewer dimensions while still preserving the similarities

More formally, we would like to compute $\cos(d, g)$ for all n docs in $\mathcal{O}(km + kn)$ ($k \ll n, m$), instead of $\mathcal{O}(nm)$

WE WILL SEE TWO ALGEBRAIC METHODS

LSI: Latent Semantic Indexing

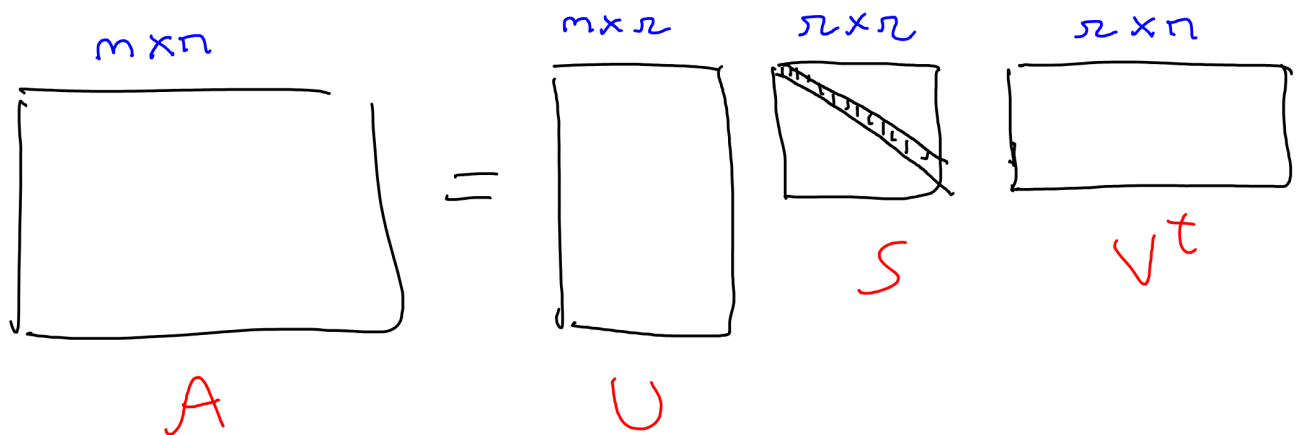
It's a data-dependent method that creates a k -dim subspace by eliminating redundant axes: this is done by "merging" related axes, like synonym and polysemy terms: therefore the axes don't represent a term anymore, but a "concept".
IMU?

In practice, it pre-processes the docs using the Singular Value Decomposition technique, and from there it creates a new smaller vector space for faster querying.

SVD

$$A = U S V^t$$

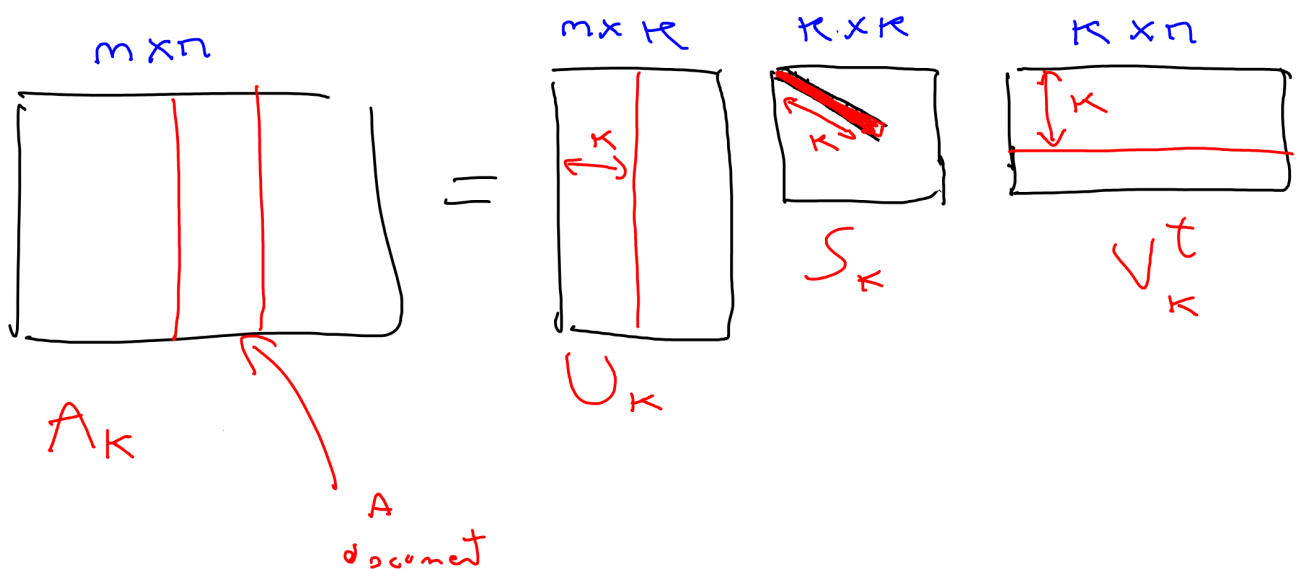
- A = $m \times n$ matrix with rank $r \leq \min(m, n)$, where row = Term and column = doc
- U = $m \times r$ matrix, where row = Term and column = eigenvector of T. T is the Term-Term correlation matrix, & symmetrical $m \times m$ matrix ($T = A A^t$)
- S = $r \times r$ diagonal matrix, that stores the eigenvalues of T in decreasing order
- V = $n \times r$ matrix, where row = doc and column = eigenvector of D. D is the doc-doc correlation matrix ($D = A^t A$)



Now let's take $\kappa \ll n$ and zero out all but the κ biggest eigenvalues in S and call S_κ this new version of S . The rank of S_κ is κ , which usually is ≈ 100 , while the rank of A is usually above 10000.

We then compute the product with U and V^T , which consider only the first κ columns of U (U_κ) and the first κ rows of V^T (V_κ^T)

$$A_\kappa = U_\kappa S_\kappa V_\kappa^T$$



A_K , of all the matrices $m \times n$ of rank K ,
is the best approximation of A , because
it preserves all the relative distances approximately

$$\|A - A_K\| = \min_{B \in \mathbb{R}^{m \times n}} \|A - B\|$$

BUT IS AS BIG AS A

Since we are interested in doc-doc similarity,
let's see how to write D now:

$$\begin{aligned} D = A^t A &= (U S V^t)^t (U S V^t) = \\ &= (S V^t)^t (S V^t) \quad n \times n \end{aligned}$$

Take $X = S V^t$ and take its rank K approximation

$$X_K = S_K V^t \quad K \times n$$

This is a good approx of D , because

$$\underline{A^t A \simeq X_K^t X_K} \quad (both \ n \times n)$$

We can use X_K to define how to project A , because

$$X_K = S_K V_K^T \Rightarrow X_K \approx U_K^T A$$

But X_K has way less rows, because $K \ll n$
therefore we've compressed the terms in "concepts"

and the columns represent the projection of the docs

Now we have to also project the query to the new space; Since we can

see the query vector as a new column

of A , $Q' = U_K^T Q$ is the projected

query (L2 norm computation)

Which are the concepts?

c -th concept = c -th col of U_K $m \times r$

• $U_K [i] [c]$ = strength of association between
 c -th concept and i -th term

• $V_K^T [c] [j]$ = strength of association
between c -th concept and j -th document

• projected document $d_j' = U_K^T d_j$

$d_j' [c]$ = strength of concept c in d_j

• projected query $q' = U_K^T q$

$q' [c]$ = strength of concept c in q

RANDOM PROJECTION

It is a data-independent reduction method that choose randomly a k -dim vector subspace that guarantees good stretching properties between any pair of points, with high probability.

It is based on the Johnson-Lindenstrauss lemma, that states that exist a function, called JL-embedding which project a vector into a smaller space preserving (almost) the distances between points.

- Let P be a set of n distinct points in m dimension ($k \ll m$)

- Given $\epsilon > 0$, exists a projection function

$f: \mathbb{R}^m \rightarrow \mathbb{R}^k$ such that:

$$\forall (u, v) \in P:$$

$$(1-\epsilon) \cdot \|u-v\|^2 \leq \|f(u) - f(v)\|^2 \leq (1+\epsilon) \|u-v\|^2$$

$$\text{And } k = O(\epsilon^{-2} \log n)$$

The lemma applies for the euclidean distance, but we can arrange it for the cosine similarity:

$$u, v \in P. \quad 2u \cdot v = \|u\|^2 + \|v\|^2 - \|u-v\|^2$$

$$2 \cdot f(u) \cdot f(v) = \underbrace{\|f(u)\|^2 + \|f(v)\|^2}_{\wedge} - \underbrace{\|f(u) - f(v)\|^2}_{\wedge \text{ LEMMA}}$$

$$2f(u) \cdot f(v) \leq (\|u\|^2 + \|v\|^2)(1+\epsilon) - \|u-v\|^2(1-\epsilon)$$

Since $\|u-v\|^2 = \|u\|^2 + \|v\|^2 - 2u \cdot v$

$$\begin{aligned} 2 \cdot f(u) \cdot f(v) &\leq (\|u\|^2 + \|v\|^2)(1+\epsilon) - (\|u\|^2 + \|v\|^2 - 2uv)(1-\epsilon) = \\ &= (\|u\|^2 + \|v\|^2)(1+\epsilon) + (2uv - \|u\|^2 - \|v\|^2)(1-\epsilon) = \\ &= (\|u\|^2 + \|v\|^2)(1+\epsilon) + 2uv - 2\epsilon uv - \|u\|^2 + \\ &\quad + \epsilon \|u\|^2 - \|v\|^2 + \epsilon \|v\|^2 = \\ &= \cancel{\|u\|^2} + \epsilon \|u\|^2 + \cancel{\|v\|^2} + \epsilon \|v\|^2 + 2uv - \\ &\quad - 2\epsilon uv - \cancel{\|u\|^2} + \epsilon \|u\|^2 - \cancel{\|v\|^2} + \epsilon \|v\|^2 = \\ &= 2\epsilon \|u\|^2 + 2\epsilon \|v\|^2 + uv(2-2\epsilon) = \\ &= \epsilon \|u\|^2 + \epsilon \|v\|^2 + uv(1-\epsilon) \end{aligned}$$

$$f(u) \cdot f(v) \leq \epsilon (\|u\|^2 + \|v\|^2) + u \cdot v (1 - \epsilon)$$

if u, v are normalized, $\|u\|, \|v\| = 1$, therefore $\cos(u, v)$ changes by at most 2ϵ

How to construct f ?

Any $n \times n$ projection matrix P w.r.t. any random distribution with mean $\mu = 0$ and variance $\sigma = 1$

EXAMPLE

$$P_{ij} = \begin{cases} 1 & \text{prob } 1/2 \\ -1 & \text{prob } 1/2 \end{cases}$$

$$P_{ij} = \begin{cases} 1 & \text{prob } 1/3 \\ -1 & \text{prob } 1/3 \\ 0 & \text{prob } 1/3 \end{cases}$$