Q1

Latent Semantic Indexing is a popular technique that uses a partial singular value decomposition (PSVD) to overcome problems of vector space representation such as synonymy – vector space representation fails to capture the relationship between synonyms (underestimating the similarity user perceives) and polysemy – vector space representation fails to represent that a term could have multiple meanings (overestimating the similarity the user perceives). We use a k-dimensional LSI representation of the matrix to overcome such problems.

Traditional LSI based approaches require us to calculate the entire PSVD again from scratch, which is a time consuming. To avoid this, incremental update of LSI is preferred.

One of the approaches is to use **folding-in**: PSVD of a term-document matrix is represented as $U_k \sum_k V_k^T$. It only updates the U_k matrix (term matrix) and the \sum_k , V_k^T (document matrix) are left unchanged. But the addition of new documents may require the update of \sum_k , V_k^T as well.

Instead, we look at a new **folding-up** approach that uses a combination of folding-in and updating at each increment in order to reduce the computational cost of updating without significantly degrading the results. The idea behind folding-up is to fold-in documents until the number of folded in documents reaches a preselected percentage of the current term-document matrix. If no updates have previously been done, the current term-document matrix is the initial matrix; otherwise it is the last updated term-document matrix. Once the number of documents that have been folded-in reaches the pre-selected percentage of the current matrix, the vectors that have been appended to V_k during folding-in are discarded. The PSVD is then updated to reflect the addition of all the document vectors that have been folded-in since the last update. These document vectors are then discarded. The process continues with folding-in until the next update.

Presenting the design and implementation of an algorithm for incrementally update of LSI:

Constructing the updated term document matrix X_2 (t+p) x (d+q) from the original term-document matrix X_1 (txd):

For terms

- The terms that are retained from the original matrix are positioned in the same row as the original
- The term that was deleted from the X_1 is still positioned in the same row, but the corresponding cells in the columns will be filled by zeros.
- A new term that is added is positioned below the rows of existing and deleted terms

For documents

- A document that is retained is in the same column as the original matrix.
- A renamed document also occupies its original position as in X₁
- A deleted document is placed in the same column as X₁ but it will be filled by zeros.
- An added document is positioned to the right of all columns of existing and deleted documents

Explaining the algorithm

The original term-document matrix X_1 can be decomposed into a product of 3 matrices using **SVD**:

$$X_1=T_{01} \cdot S_{01} \cdot D_{01}$$
 (1)

 S_{01} is the diagonal matrix of singular values. T_{01} is the matrix of eigenvectors of the square symmetric matrix $X_1 \cdot X_1'$ (X_1 ' is the transpose of X_1). T_{01} is unitary. Since $X_1 \cdot X_1'$ is real, T_{01} has orthonormal columns and rows, i.e., $T_{01}' \cdot T_{01} = I$ and $T_{01} \cdot T_{01}' = I$. Meanwhile, D_{01} is the matrix of eigenvectors of $X_1' \cdot X_1$. Similarly, D01 also has orthonormal columns and rows, i.e., $D_{01}' \cdot D_{01} = I$ and $D_{01} \cdot D_{01}' = I$. Note that T_{01} is a $t \times m$ matrix, S_{01} is a $t \times m$ matrix, and $T_{01}' \cdot T_{01} = I$ and $T_{01}' \cdot T_{01} = I$ and $T_{01}' \cdot T_{01} = I$.

If the singular values in S_{01} can be ordered by size, the first k largest may be kept and the remaining smaller ones set to zero. Since zeroes were introduced into S_{01} , the representation can be simplified by deleting the zero rows and columns of S_{01} to obtain a new diagonal matrix S_1 , and then by deleting the corresponding columns of T_{01} and D_{01} to obtain T_1 and T_2 such that

$$X_1 \approx X_1^a = T_1 \cdot S_1 \cdot D_1'$$
 (2)

X, = X, a = T, S, D,

T, is t x k

S, is k x k

D, is k x k

D, is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

X, a is matrix of rank k, which closest to X, (best pass; ble

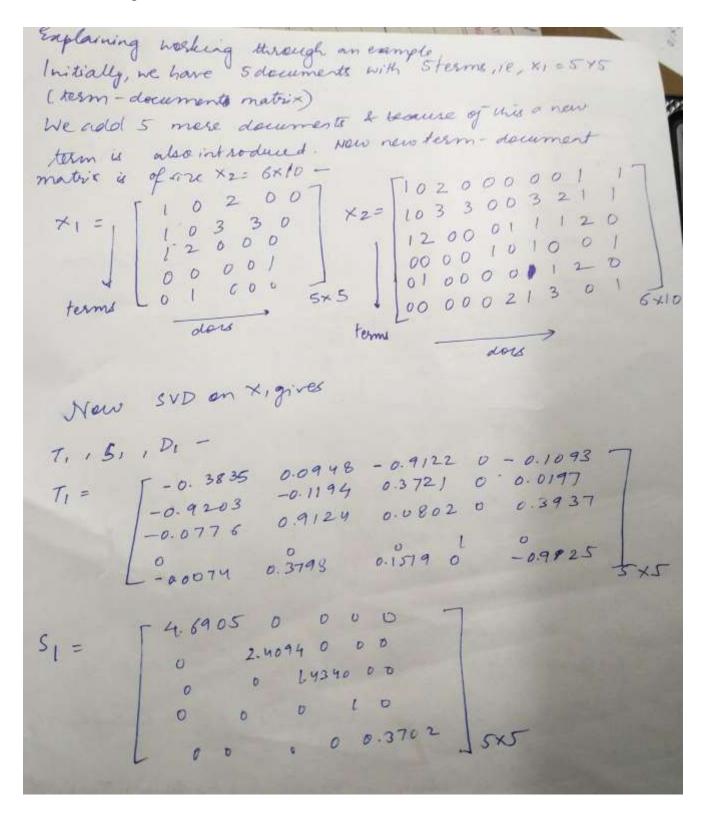
X, a is matrix of rank k, which cl

X1 - original term document matrix (txd) X2 - new term document matrix (++1) x (d+q) his want X, to have same size as X2, by transforming X, into an expanded matrix XIE of dimension (+p)x(d+q) To achieve this, he take I matrices V of size (t+p) + t and Vand V are extensions of Identity matrices It and Identity with prous & a columns of all recos, respectively. V of size dr(d+2). New, to construct XIE = U.XI.V But from (1) this can be written as-XIE= U.XI.V= U. To1. Soi. Doi. V = (U.Tor). (Sor VDa) (U.To1). So1. (V.Do1) XIE = Tole . Sol. D'ore where Tore = D. Tor, Dore = V' Tole. Tole = (U.Ta). (U.Ta) = U.Tal. Tol. U' = U.Z.U'=U.U'_I Comitably Toil. Toile = I, Doile - Doile = I, D'ale . Doile = I => We can compute an SVD for X1e from SVD of XIE, nevan perfer m reduction by heeping to largest singular value as in XI. Xie = Xie = Tie. Si. Die -(4) : SVD of XIE 7 XIE is obtained. XIE & Xihave same size of (t+p) x (d+9) : Xz = XietA where A is modification matrix

We can represent X2 = X1e + A as X2 & X1e are of Same dimensions. - (5) P= Tie. ×2. Die Tie -> k x (t+p) matrix X2 -> (+p) x (d+q) matrix Die -> (a+q) + k matrix P's dimension is kxk. Applying SVD on P, neget P= Tie. x2. Die=Tung. Sz. D'emp Tie. Tie. X2. Die. Die = Tie. Was. Trong. Sz. D'rong Die => X2= Tie Ttemp. Sz. Dtemp. Die - (6) => | X2 = To2. S2. Do2 where To2 = Tie Trans Doz = Die Dramp Incremental Update of LSI -Obtain X2 as described (2) Compute Tole = U. Tol 4 Dole = 11. Duive Tie, Die (3) Compute P= Tie. X2. Die. Perform SVD to get P= Treng. Sz. Dremp (Pick+k) (5) Compute To2, Do2' Derive T2 + D2 by the reduchas step Thue, we have found out a new decomposition T2, 52, 0's for X2 x X2.

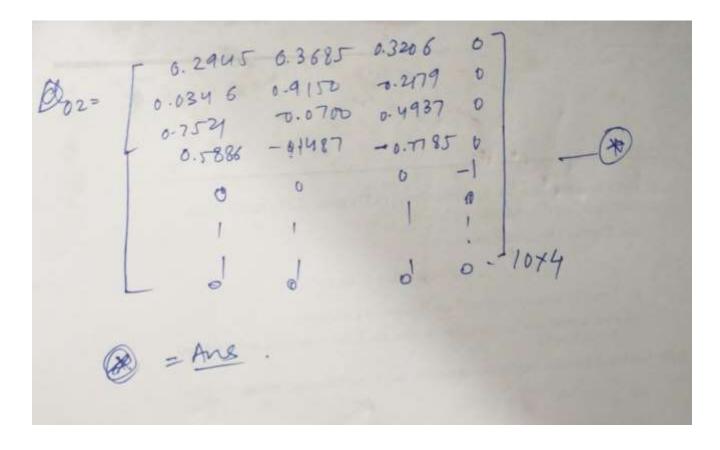
This is the incremental update of the LSI that saves time by not requiring one to reconstruct the SVD from scratch.

Explaining the working of the algorithm through an example of 5 documents, showing how adding 5 more documents change the document and term vectors. [Code is attached in folder for calculations shown]:



```
-0.2945 -0.3685 0.3206 0 0.8214
Doie =
          -0.0346
                  -0-9150
                          -0.2179 0 -0.3378
                                    -0.4301
          -0.7521
                  0-0700
                          0.4937
          -0.5886 0-1487 -0.77850 0-1596
                                    0
                    10
            0
                                               0
                         0
             6
             0
                                                     10--
                                                        10410
                         0 0 6x4
                  2.4094 00
                 0 1.434 0
                     -0.3685 0.3206
            -0.2945
                             -0.2179
                      -0.9150
  Die=
            -0.346
                               1.4437
                     0.0700
            -0-7521
                     0.1487
            -0-1886
                              0.9122
                             -0-374 0
            -0.3835
                              -0.0802 0
             -0.4203
   Tie=
                      -0.9124
                              0
              -0.0776
                              -0.1579
                       -0.5748
              -0.074
                              - 4.44c-16
                              -4.1633E-16
           4.6905
    D=
                    240948
                                 1.4340
                    -8326 e+6
            2 2204 e-16
                                            4×4
```

New SVD of P Tremp S2 Premp Tranp= [-2.958e-17 -4.7339e-17]
-2.95e-17 -4 6.903e-16 |
-4.9339e-17 6.903e-16 | 0 S2= T 0 0 0 0 -0 + 5-54e+6 00 - 5-551e+6 1 0 0.3835 0.0948 6.9122 1 0.9203 -0.1194 -0.3724 0 0.0776 0.9124 -0.0802 0 0.0074 0.9124 0.0802 0 To2= 0



Q2Various Clustering techniques can be useful in Information Retrieval Tasks:

- Search Result Clustering One of the most useful applications in Information Retrieval based systems is to cluster the search results of user query. For example, If the user searches for 'apple' on Google it should not only display the fruit apple but the Apple brand and cluster similar results together. This will speed up the lookup and user will be able to identify the required information.
- Scatter-Gather Clusters the whole collection to get groups of documents that the user can select or gather. These selected groups are further merged and result is clustered. Result is repeated until a cluster of interest is found. This helps in searching without typing by providing alternative user interface.
- Cluster-based navigation Many a times, users prefer browsing over searching when they are
 unsure of which search terms to use. This is an alternative to user-mediated scatter-gather technique
 that requires computing a static hierarchical clustering of a collection that is not influenced by user
 interactions.
- Cluster hypothesis to improve search results Using standard inverted index to identify an initial set of documents that match the query, that are merged with other documents from the same cluster. Example, car matches to automobile, vehicle. So, apart from taking documents containing the term car, we use documents that contain the related words like vehicle, automobile.
- K-means based clustering for improved web search User's query goes to a meta search engine and retrieves the relevant documents. It also uses the terms in the query to find their co-occuring terms from Wikipedia or related sites, and compares each co-occurent term in the document. The count of

co-occurent terms helps to weight the document. Next, a knowledge base for clustering is formed as a meta-directory tree with the obtained search results. Using the meta-directory tree, clustering is performed with K-Means. We initialize the number of clusters and assign thresholds for each. Then, calculate the weight of a web-page using outlinks and inlinks in the page. Use a threshold to assign the webpages to clusters, and repeat until all pages are clustered. Finally, the cluster with highest threshold is returned. If a node has more number of relevant links, it is treated as the centroid of the cluster. The centroid value is determined by considering the weight of a document.

• Language Modelling – Language modelling is a probability distribution over sequence of words. Given such a sequence, it assigns a probability to the whole sequence. Clustering helps to avoid sparse data problems in the language modeling approach. The model of a document d can be interpolated with a collection model. If the collection model contains many documents with terms not in d, the collection model can be replaced by a model derived from d's cluster and get an accurate representation of the probability distribution of terms in d.

These techniques can help improve performance over the traditional IR tasks based on inverted index:

- In standard IR systems, search results are presented in a simple list that user has to scroll down all
 the way to find the required result. Search result clustering clusters the results so that similar
 documents appear together.
- Instead of searching for keywords, cluster based navigation is useful in scenarios where they are unsure of what search terms to use and hence prefer browsing over searching.
- Clustering is well suited for access to a collection of news stories because news stories is not really search, but rather a process of selecting a subset of stories about recent events.
- Clustering helps to speed up search. Inverted index supports fast nearest-neighbor search for the standard IR setting. But we may not be able to use this inverted index efficiently, like in latent semantic indexing synonymy: we cannot identify the documents containing synonyms to the query term which may be useful in providing relevant information; polysemy: we cannot identify the context/ meaning in which the term is used in the document. Using cluster hypothesis, we can find the clusters closest to the query and only consider the documents from these clusters. Within this smaller set, we can compute similarity and rank documents in the usual way. There are fewer clusters than the documents, making the process of finding the closest cluster fast. The documents matching the query would be in the same cluster, reducing the search time over traditional inverted index.