

LEMP: Fast Retrieval of Large Entries in a Matrix Product

Christina Teflioudi¹ Rainer Gemulla² Olga Mykytiuk

¹Max Planck Institute for Computer Science
Saarbrücken, Germany

²Mannheim University
Mannheim, Germany

May 4, 2015



MAX-PLANCK-GESELLSCHAFT

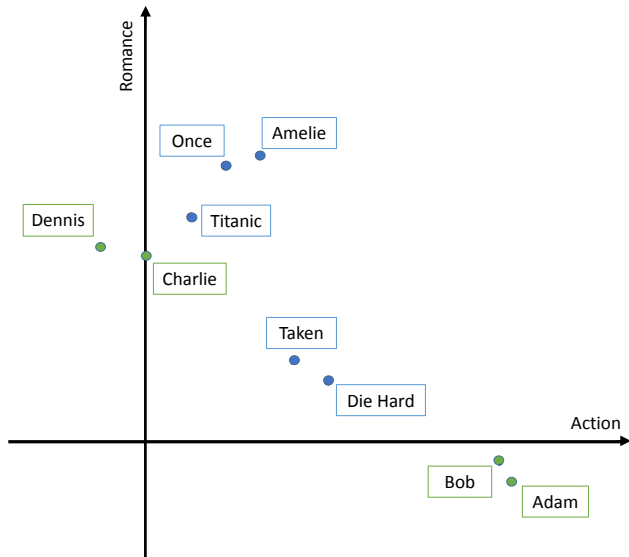
mpii max planck institut
informatik

Recommender Systems

- ▶ Problem
 - ▶ Set of users
 - ▶ Set of items (movies, books, jokes, products, stories, ...)
 - ▶ Feedback (ratings, purchase, click-through, tags, ...)
- ▶ Predict the preference of each user for each item

	Die Hard	Taken	Once	Amelie	Titanic
Adam	5	?	1	2	?
Bob	5	4	?	?	1
Charlie	2	?	5	?	4
Dennis	?	1	5	5	?

Latent Factor Models



Recommender Systems

- ▶ Given
 - ▶ A query matrix \mathbf{Q}

$$\begin{array}{l} \textit{Adam} \\ \textit{Bob} \\ \textit{Charlie} \\ \textit{Dennis} \end{array} \begin{pmatrix} 3.2 & -0.4 \\ 3.1 & -0.2 \\ 0 & 1.8 \\ -0.4 & 1.9 \end{pmatrix} \\ \mathbf{Q}^T$$

Recommender Systems

- ▶ Given
 - ▶ A query matrix \mathbf{Q}
 - ▶ A probe matrix \mathbf{P}

$$\begin{array}{l} \textit{Adam} \\ \textit{Bob} \\ \textit{Charlie} \\ \textit{Dennis} \end{array} \begin{pmatrix} 3.2 & -0.4 \\ 3.1 & -0.2 \\ 0 & 1.8 \\ -0.4 & 1.9 \end{pmatrix} \mathbf{Q}^T$$

$$\begin{array}{ccccc} & \textit{Die Hard} & \textit{Taken} & \textit{Once} & \textit{Amelie} & \textit{Titanic} \\ \begin{pmatrix} 1.6 & 1.3 & 0.7 & 1 & 0.4 \\ 0.6 & 0.8 & 2.7 & 2.8 & 2.2 \end{pmatrix} & & & & & \end{array} \mathbf{P}$$

Recommender Systems

- ▶ Given
 - ▶ A query matrix \mathbf{Q}
 - ▶ A probe matrix \mathbf{P}
 - ▶ A threshold $\theta > 0$
- ▶ Find good recommendations

$$\begin{array}{l} \textit{Adam} \\ \textit{Bob} \\ \textit{Charlie} \\ \textit{Dennis} \end{array} \begin{pmatrix} 3.2 & -0.4 \\ 3.1 & -0.2 \\ 0 & 1.8 \\ -0.4 & 1.9 \end{pmatrix} \mathbf{Q}^T$$

$$\begin{array}{ccccc} & \textit{Die Hard} & \textit{Taken} & \textit{Once} & \textit{Amelie} & \textit{Titanic} \\ \begin{pmatrix} 1.6 & 1.3 & 0.7 & 1 & 0.4 \\ 0.6 & 0.8 & 2.7 & 2.8 & 2.2 \end{pmatrix} & & & & & \end{array} \mathbf{P}$$

Recommender Systems

- ▶ Given
 - ▶ A query matrix \mathbf{Q}
 - ▶ A probe matrix \mathbf{P}
 - ▶ A threshold $\theta > 0$
- ▶ Find good recommendations
 - ▶ All entries in $\mathbf{Q}^T \mathbf{P}$ that are $\geq \theta$

		Die Hard	Taken	Once	Amelie	Titanic	
		1.6	1.3	0.7	1	0.4	\mathbf{P}
		0.6	0.8	2.7	2.8	2.2	
Adam	\mathbf{Q}^T	4.9	3.8	1.2	2.1	0.4	$\mathbf{Q}^T \mathbf{P}$
Bob		4.8	3.9	1.6	2.5	0.8	
Charlie		1	1.4	4.9	5.0	4.0	
Dennis		0.5	1	4.9	4.9	4.0	

Recommender Systems

- ▶ Given
 - ▶ A query matrix \mathbf{Q}
 - ▶ A probe matrix \mathbf{P}
 - ▶ A threshold $\theta > 0$
- ▶ Find good recommendations
 - ▶ All entries in $\mathbf{Q}^T \mathbf{P}$ that are $\geq \theta$
 - ▶ Each entry is an inner product $\mathbf{q}^T \mathbf{p} = \sum_{i=1}^r q_i p_i$

		Die Hard	Taken	Once	Amelie	Titanic	
		1.6	1.3	0.7	1	0.4	\mathbf{P}
		0.6	0.8	2.7	2.8	2.2	
Adam	\mathbf{Q}^T	4.9	3.8	1.2	2.1	0.4	$\mathbf{Q}^T \mathbf{P}$
Bob		4.8	3.9	1.6	2.5	0.8	
Charlie		1	1.4	4.9	5.0	4.0	
Dennis		0.5	1	4.9	4.9	4.0	

Problem Statement

Maximum Inner Product Search

Find pairs of vectors with large inner products

Given

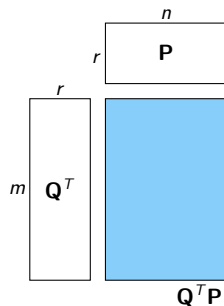
- ▶ a **query vector** \mathbf{q}^T (from matrix \mathbf{Q}^T)
- ▶ a set \mathbf{P} of **probe vectors** (the matrix \mathbf{P})
- ▶ a threshold $\theta > 0$

Find

- ▶ all vectors \mathbf{p} such that $\mathbf{q}^T \mathbf{p} \geq \theta$

Naive Solution

- ▶ Compute full matrix product $\mathbf{Q}^T \mathbf{P}$
- ▶ Determine which entries are $\geq \theta$
- ▶ Complexity $O(mnr)$
- ▶ Usually
 - ▶ m, n : order of millions
 - ▶ $10 < r < 500$
- ▶ Example
 - ▶ $m = 10$ millions
 - ▶ $n = 1$ million
 - ▶ #entries = 10 trillion
 - ▶ Avg. Inner Product Time = 100nsec
 - ▶ Runtime > 11 days
- ▶ Can we do better than that?



When is an inner product large?

$$\mathbf{q}^T \mathbf{p} = \|\mathbf{q}\| \|\mathbf{p}\| \cos \angle(\mathbf{q}, \mathbf{p}) \geq \theta, \quad -1 \leq \cos \angle(\mathbf{q}, \mathbf{p}) \leq 1$$

When is an inner product large?

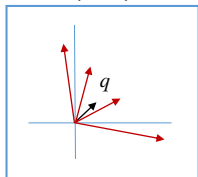
$$\mathbf{q}^T \mathbf{p} = \|\mathbf{q}\| \|\mathbf{p}\| \cos \angle(\mathbf{q}, \mathbf{p}) \geq \theta, \quad -1 \leq \cos \angle(\mathbf{q}, \mathbf{p}) \leq 1$$

Length
($\|\mathbf{q}\| \|\mathbf{p}\|$)

short
($< \theta$)

medium
($\approx \theta$)

long
($\gg \theta$)



Angle
($\cos \angle(\mathbf{q}, \mathbf{p})$)

small positive
large positive

Suitable
method

When is an inner product large?

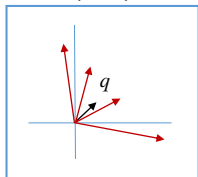
$$\mathbf{q}^T \mathbf{p} = \|\mathbf{q}\| \|\mathbf{p}\| \cos \angle(\mathbf{q}, \mathbf{p}) \geq \theta, \quad -1 \leq \cos \angle(\mathbf{q}, \mathbf{p}) \leq 1$$

Length
($\|\mathbf{q}\| \|\mathbf{p}\|$)

short
($< \theta$)

medium
($\approx \theta$)

long
($\gg \theta$)



Angle
($\cos \angle(\mathbf{q}, \mathbf{p})$)

small positive
large positive

Suitable
method

Prune

When is an inner product large?

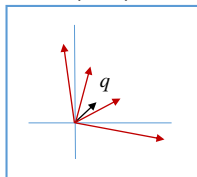
$$\mathbf{q}^T \mathbf{p} = \|\mathbf{q}\| \|\mathbf{p}\| \cos \angle(\mathbf{q}, \mathbf{p}) \geq \theta, \quad -1 \leq \cos \angle(\mathbf{q}, \mathbf{p}) \leq 1$$

Length
($\|\mathbf{q}\| \|\mathbf{p}\|$)

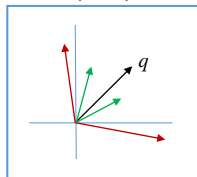
short
($< \theta$)

medium
($\approx \theta$)

long
($\gg \theta$)



small positive
large positive



small positive
large positive

Angle
($\cos \angle(\mathbf{q}, \mathbf{p})$)

Suitable
method

Prune

When is an inner product large?

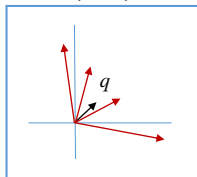
$$\mathbf{q}^T \mathbf{p} = \|\mathbf{q}\| \|\mathbf{p}\| \cos \angle(\mathbf{q}, \mathbf{p}) \geq \theta, \quad -1 \leq \cos \angle(\mathbf{q}, \mathbf{p}) \leq 1$$

Length
($\|\mathbf{q}\| \|\mathbf{p}\|$)

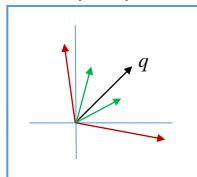
short
($< \theta$)

medium
($\approx \theta$)

long
($\gg \theta$)



small positive
large positive



small positive
large positive

Angle
($\cos \angle(\mathbf{q}, \mathbf{p})$)

Suitable
method

Prune

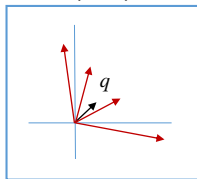
Cosine Similarity
Search $\theta' = \frac{\theta}{\|\mathbf{q}\| \|\mathbf{p}\|}$

When is an inner product large?

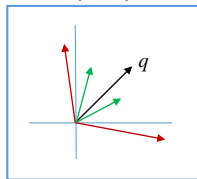
$$\mathbf{q}^T \mathbf{p} = \|\mathbf{q}\| \|\mathbf{p}\| \cos \angle(\mathbf{q}, \mathbf{p}) \geq \theta, \quad -1 \leq \cos \angle(\mathbf{q}, \mathbf{p}) \leq 1$$

Length
($\|\mathbf{q}\| \|\mathbf{p}\|$)

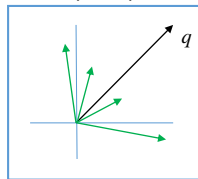
short
($< \theta$)



medium
($\approx \theta$)



long
($\gg \theta$)



Angle
($\cos \angle(\mathbf{q}, \mathbf{p})$)

small positive
large positive

small positive
large positive

small positive
large positive

Suitable
method

Prune

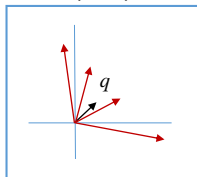
Cosine Similarity
Search $\theta' = \frac{\theta}{\|\mathbf{q}\| \|\mathbf{p}\|}$

When is an inner product large?

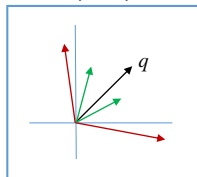
$$\mathbf{q}^T \mathbf{p} = \|\mathbf{q}\| \|\mathbf{p}\| \cos \angle(\mathbf{q}, \mathbf{p}) \geq \theta, \quad -1 \leq \cos \angle(\mathbf{q}, \mathbf{p}) \leq 1$$

Length
($\|\mathbf{q}\| \|\mathbf{p}\|$)

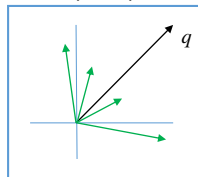
short
($< \theta$)



medium
($\approx \theta$)



long
($\gg \theta$)



Angle
($\cos \angle(\mathbf{q}, \mathbf{p})$)

small positive
large positive

small positive
large positive

small positive
large positive

Suitable
method

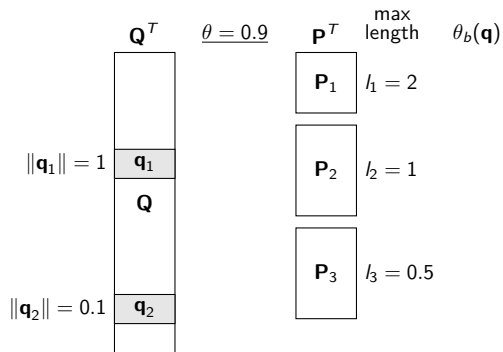
Prune

Cosine Similarity
Search $\theta' = \frac{\theta}{\|\mathbf{q}\| \|\mathbf{p}\|}$

Naive-like
retrieval

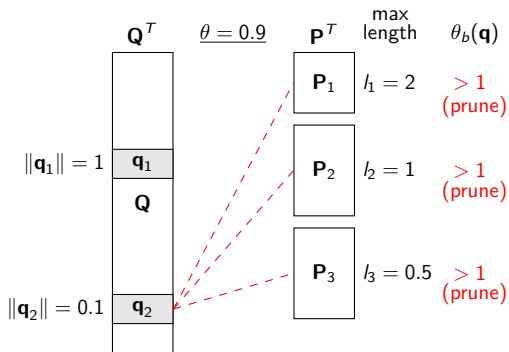
Main idea: bucketize by length

- ▶ Partition \mathbf{P} in buckets with vectors of similar length
- ▶ Index buckets suitably



Main idea: bucketize by length

- ▶ Partition \mathbf{P} in buckets with vectors of similar length
- ▶ Index buckets suitably
- ▶ For each query vector and bucket
 - ▶ Determine local threshold
 - ▶ Prune bucket if possible

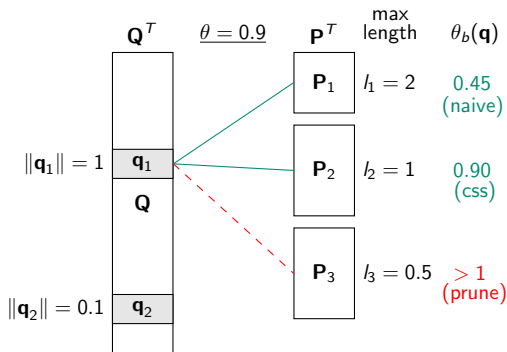


Local threshold on cosine similarity:

$$\theta_b(\mathbf{q}) = \frac{\theta}{\|\mathbf{q}\|l_b}$$

Main idea: bucketize by length

- ▶ Partition \mathbf{P} in buckets with vectors of similar length
- ▶ Index buckets suitably
- ▶ For each query vector and bucket
 - ▶ Determine local threshold
 - ▶ Prune bucket if possible
 - ▶ Otherwise, select best retrieval method



Local threshold on cosine similarity:

$$\theta_b(\mathbf{q}) = \frac{\theta}{\|\mathbf{q}\|l_b}$$

Discussion

LEMP

- ▶ Prunes whole buckets
- ▶ Selects a suitable retrieval algorithm per query and bucket
- ▶ Can leverage existing methods
- ▶ Cache-friendly

Bucket-level retrieval

Choose among a variety of algorithms

- ▶ Fagin's Threshold Algorithm ($r < 10$)
- ▶ All pairs similarity search family ($1000 < r$)
- ▶ Space-partitioning trees ($r < 10$)

Our vectors are

- ▶ Not necessarily sparse
- ▶ Real values
- ▶ Medium dimensionality ($10 < r < 500$)

Two new algorithms

- ▶ COORD
- ▶ INCR

INCR: Main Idea

- ▶ \mathbf{q}, \mathbf{p} qualify if $\|\mathbf{q}\| \|\mathbf{p}\| (\sum_{i=1}^r \bar{q}_i \bar{p}_i) \geq \theta$

$\ \mathbf{q}\ $	\bar{q}_1	\bar{q}_2	\bar{q}_3	\bar{q}_4	\dots	\bar{q}_r
$\ \mathbf{p}\ $	\bar{p}_1	\bar{p}_2	\bar{p}_3	\bar{p}_4	\dots	\bar{p}_r

INCR: Main Idea

- ▶ \mathbf{q}, \mathbf{p} qualify if $\|\mathbf{q}\| \|\mathbf{p}\| (\sum_{i=1}^r \bar{q}_i \bar{p}_i) \geq \theta$
- ▶ Assume you have a budget of $\phi = 3$ multiplications for the $\sum_{i=1}^r \bar{q}_i \bar{p}_i$

$\ \mathbf{q}\ $	\bar{q}_1	\bar{q}_2	\bar{q}_3	\bar{q}_4	\dots	\bar{q}_r
$\ \mathbf{p}\ $	\bar{p}_1	\bar{p}_2	\bar{p}_3	\bar{p}_4	\dots	\bar{p}_r

INCR: Main Idea

- ▶ \mathbf{q}, \mathbf{p} qualify if $\|\mathbf{q}\|\|\mathbf{p}\|(\sum_{i=1}^r \bar{q}_i \bar{p}_i) \geq \theta$
- ▶ Assume you have a budget of $\phi = 3$ multiplications for the $\sum_{i=1}^r \bar{q}_i \bar{p}_i$
- ▶ Goal: decide after seeing $\phi = 3$ coordinates:

$\ \mathbf{q}\ $	\bar{q}_1	\bar{q}_2	\bar{q}_3	\bar{q}_4	\dots	\bar{q}_r
$\ \mathbf{p}\ $	\bar{p}_1	\bar{p}_2	\bar{p}_3	\bar{p}_4	\dots	\bar{p}_r

INCR: Main Idea

- ▶ \mathbf{q}, \mathbf{p} qualify if $\|\mathbf{q}\|\|\mathbf{p}\|(\sum_{i=1}^r \bar{q}_i \bar{p}_i) \geq \theta$
- ▶ Assume you have a budget of $\phi = 3$ multiplications for the $\sum_{i=1}^r \bar{q}_i \bar{p}_i$
- ▶ Goal: decide after seeing $\phi = 3$ coordinates:
 $\|\mathbf{q}\|\|\mathbf{p}\|(\sum_{i=1}^3 \bar{q}_i \bar{p}_i + \text{UpperBoundFor}(\sum_{i=4}^r \bar{q}_i \bar{p}_i)) \geq \theta$

$\ \mathbf{q}\ $	\bar{q}_1	\bar{q}_2	\bar{q}_3	\bar{q}_4	...	\bar{q}_r
$\ \mathbf{p}\ $	\bar{p}_1	\bar{p}_2	\bar{p}_3	\bar{p}_4	...	\bar{p}_r

INCR: Main Idea

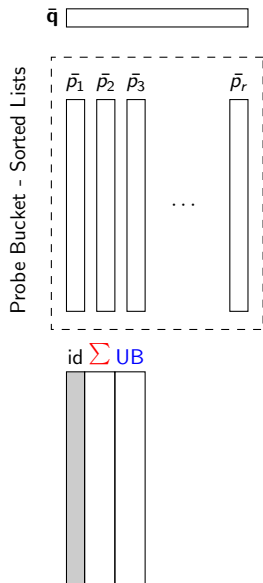
- ▶ \mathbf{q}, \mathbf{p} qualify if $\|\mathbf{q}\|\|\mathbf{p}\|(\sum_{i=1}^r \bar{q}_i \bar{p}_i) \geq \theta$
- ▶ Assume you have a budget of $\phi = 3$ multiplications for the $\sum_{i=1}^r \bar{q}_i \bar{p}_i$
- ▶ Goal: decide after seeing $\phi = 3$ coordinates:
 $\|\mathbf{q}\|\|\mathbf{p}\|(\sum_{i=1}^3 \bar{q}_i \bar{p}_i + \text{UpperBoundFor}(\sum_{i=4}^r \bar{q}_i \bar{p}_i)) \geq \theta$
- ▶ In practice
 - ▶ ϕ is automatically tuned
 - ▶ the ϕ coordinates do not have to be consecutive

$\ \mathbf{q}\ $	\bar{q}_1	\bar{q}_2	\bar{q}_3	\bar{q}_4	...	\bar{q}_r
$\ \mathbf{p}\ $	\bar{p}_1	\bar{p}_2	\bar{p}_3	\bar{p}_4	...	\bar{p}_r

INCR

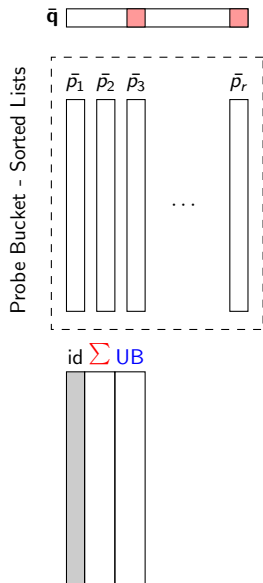
\bar{q}

INCR



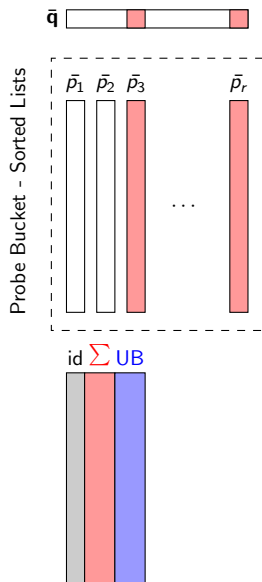
INCR

- ▶ Pick coordinates (largest $|\bar{q}_i|$ first)



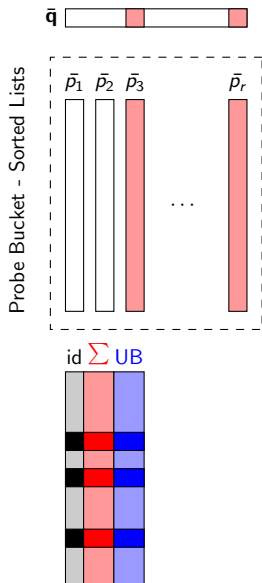
INCR

- ▶ Pick coordinates (largest $|\bar{q}_i|$ first)
- ▶ Scan indexes and update “ Σ ” and “UB” quantities



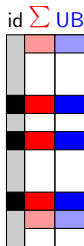
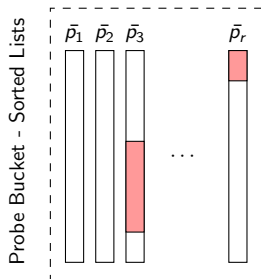
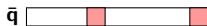
INCR

- ▶ Pick coordinates (largest $|\bar{q}_i|$ first)
- ▶ Scan indexes and update “ Σ ” and “UB” quantities
- ▶ Prune vectors for which $\|\mathbf{q}\| \|\mathbf{p}\| (\Sigma + UB) < \theta$



INCR

- ▶ Pick coordinates (largest $|\bar{q}_i|$ first)
- ▶ Scan indexes and update “ Σ ” and “UB” quantities
- ▶ Prune vectors for which $\|\mathbf{q}\| \|\mathbf{p}\| (\Sigma + UB) < \theta$
- ▶ No need to scan the whole lists!
Bounds exist

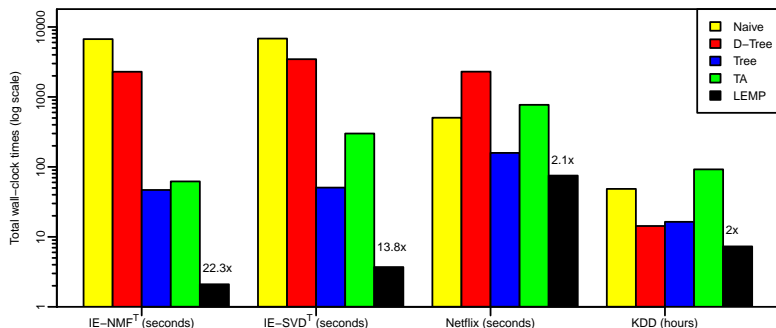


Discussion

INCR

- ▶ Computes partial inner products
- ▶ Uses simple index - cheap to construct
- ▶ Scans part of the index - does not necessarily scan from top
- ▶ Sequential memory access pattern, fast

How fast is it? (Top-1 movie per user)



rows
cols
Sparse?
Length skew
Dimensions

132K
771K
Yes
Strong
50

132K
771K
No
Strong
50

480K
17K
No
Mild
50

1000K
624K
No
Mild
50

Summary

Large inner-product search

- ▶ Matrix factorization common technique in data mining
- ▶ Large entries in matrix products are usually of particular interest

The LEMP algorithm

- ▶ Bucketizes vectors by length
- ▶ Prunes buckets whenever possible
- ▶ For the remaining buckets: selects efficient retrieval algorithms
- ▶ Consistently fastest bucket method: INCR

Summary

Large inner-product search

- ▶ Matrix factorization common technique in data mining
- ▶ Large entries in matrix products are usually of particular interest

The LEMP algorithm

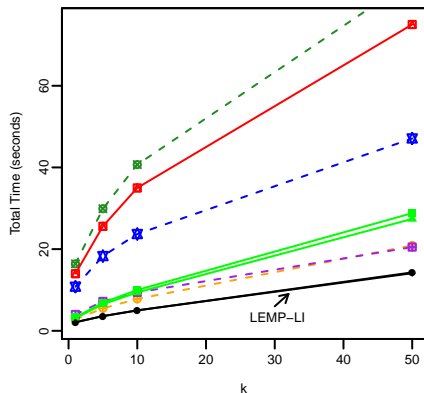
- ▶ Bucketizes vectors by length
- ▶ Prunes buckets whenever possible
- ▶ For the remaining buckets: selects efficient retrieval algorithms
- ▶ Consistently fastest bucket method: INCR

Thank you!

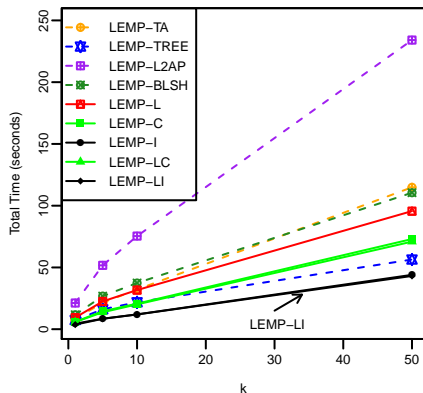
Questions?

Performance of bucket algorithms

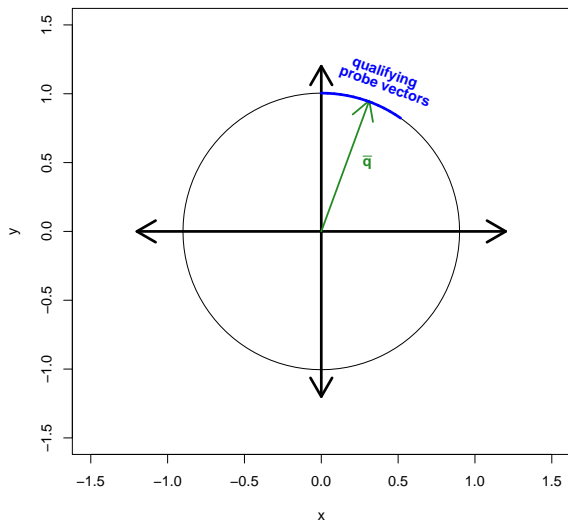
IE-NMF^T



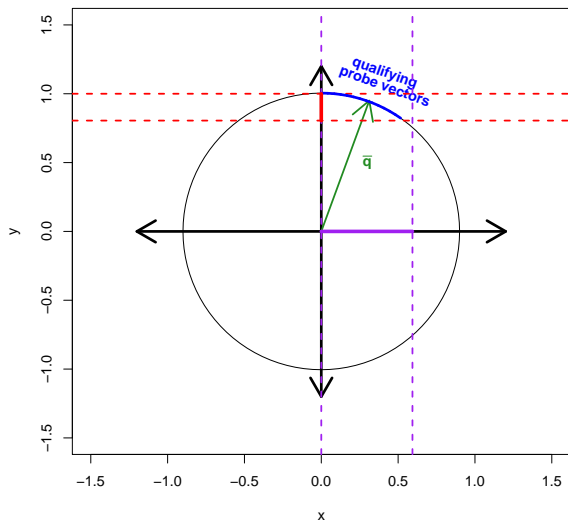
IE-SVD^T



Do we need to scan the whole lists?

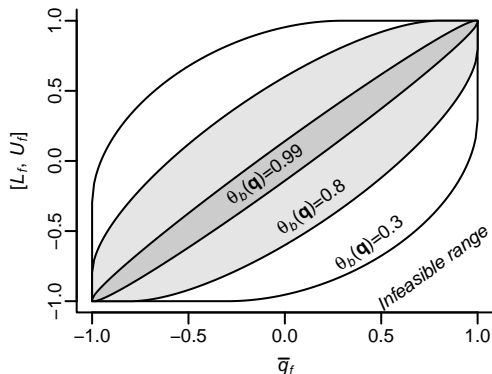


Do we need to scan the whole lists?



COORD: bounds

- Find for each coordinate \bar{p}_f of $\bar{\mathbf{p}}$ an upper and lower bound $[L_f, U_f]$ such that $\bar{p}_f \notin [L_f, U_f] \Rightarrow \bar{\mathbf{q}}^T \bar{\mathbf{p}} < \theta_b(\mathbf{q})$



Algorithm selection

- ▶ Take a sample of queries
- ▶ Run both naive and INCR/COORD
- ▶ Estimate the value of t_b
- ▶ Linear classifier with $\theta_b(\mathbf{q})$ as feature

