

*SKIP-AND-PRUNE:  
Cosine-based Top-K  
Query Processing  
for Efficient Context-Sensitive  
Document Retrieval*

Jong Wook Kim

K. Selçuk Candan

# Obiettivi in teoria...

2

Modello dei dati  
basato su parole  
chiave



Contesto di ricerca  
dell'utente



Algoritmo in grado di utilizzare questi  
modelli per restituire all'utente i  
migliori documenti che soddisfano i  
parametri di ricerca

# ...e in pratica

3

Facciamo un esempio...

Mettiamoci nei panni di un appassionato di musica che voglia ricercare maggiori informazioni

VIOLA

Enter  
←



# Obiettivo dell'algoritmo

4

- Utilizzare abilmente tecniche top-k per **reinterpretare** a seconda del contesto utente, solo i documenti necessari e limitare l'esplorazione di grandi porzioni del database
- Uso di una scoring function basata sul coseno quindi non monotona, per comparare i documenti del db con la query dell'utente



**necessità di un nuovo algoritmo  
di query processing**



$$sim_{cos}(\vec{d}, \vec{q}) = \cos(\vec{d}, \vec{q}) = \frac{\vec{d} \cdot \vec{q}}{\|\vec{d}\| \|\vec{q}\|}$$

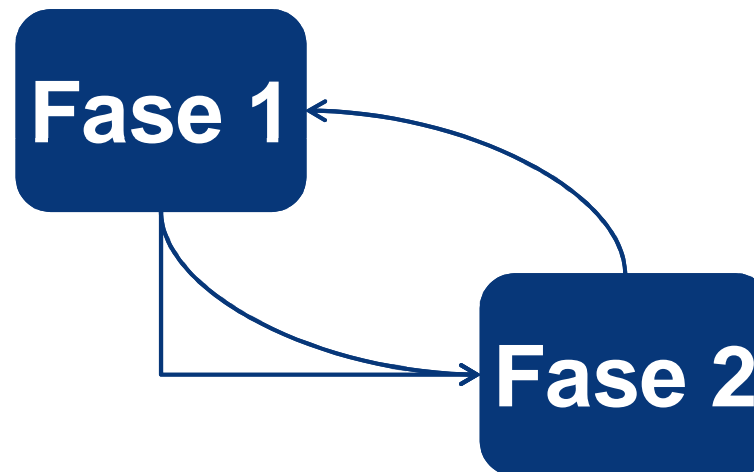


$$sim'_{cos}(\vec{d}, \vec{q}, U) = \cos(U\vec{d}, U\vec{q})$$

# Skip-and-Prune (SnP)

5

- L'algoritmo opera in due fasi:
  - ▣ FASE 1: selezione dei documenti secondo lo specifico contesto
  - ▣ FASE 2: scelta dei migliori k documenti in base alla query



- La soluzione al problema consiste nel trovare i k documenti in DB, con il più alto punteggio di similarità del coseno rispetto alla query considerando lo spazio specifico dei concetti dell'utente

# Formalizzazione del problema

6

- $\mathcal{L} \rightarrow$  dizionario contenente  $l$  distinte parole chiave

ID_KEYWORD	KEYWORD
I1	“viola”
I2	“spartito”
I3	“inno”
I4	“cantare”
I5	“calcio”



# Formalizzazione del problema

7

- $DB \rightarrow$  collezione di  $m$  documenti, rappresentati con la matrice  $D = m \times l$



D	l1	l2	l3	l4	l5
concerto	0,30	0,40	0,10	0,10	0,10
flauto_magico_mozart	0,20	0,38	0,09	0,21	0,00
musica_nello_sport	0,10	0,11	0,30	0,20	0,20
accordare_viola	0,81	0,13	0,06	0,00	0,00
fiorentina_firenze	0,29	0,00	0,08	0,00	0,40

- $D$  quindi rappresenta quanto sia realmente presente una parola chiave in un certo documento

# Formalizzazione del problema

8

- $U \rightarrow$  matrice del contesto dell'utente che riassume quanto la parola chiave sia importante all'interno di un determinato concetto per lo specifico utente

U	I1	I2	I3	I4	I5
musica	0,28	0,30	0,12	0,30	0,00
sport	0,10	0,00	0,20	0,00	0,70



- Il risultato del prodotto  $D_u = DU^T$  è la matrice rappresentante quanto un dato documento sia importante all'interno di un concetto in un contesto preciso
- Molto **oneroso**, le dimensioni delle matrici rendono incalcolabile
- Qui entrano in gioco le due fasi dell'algoritmo!

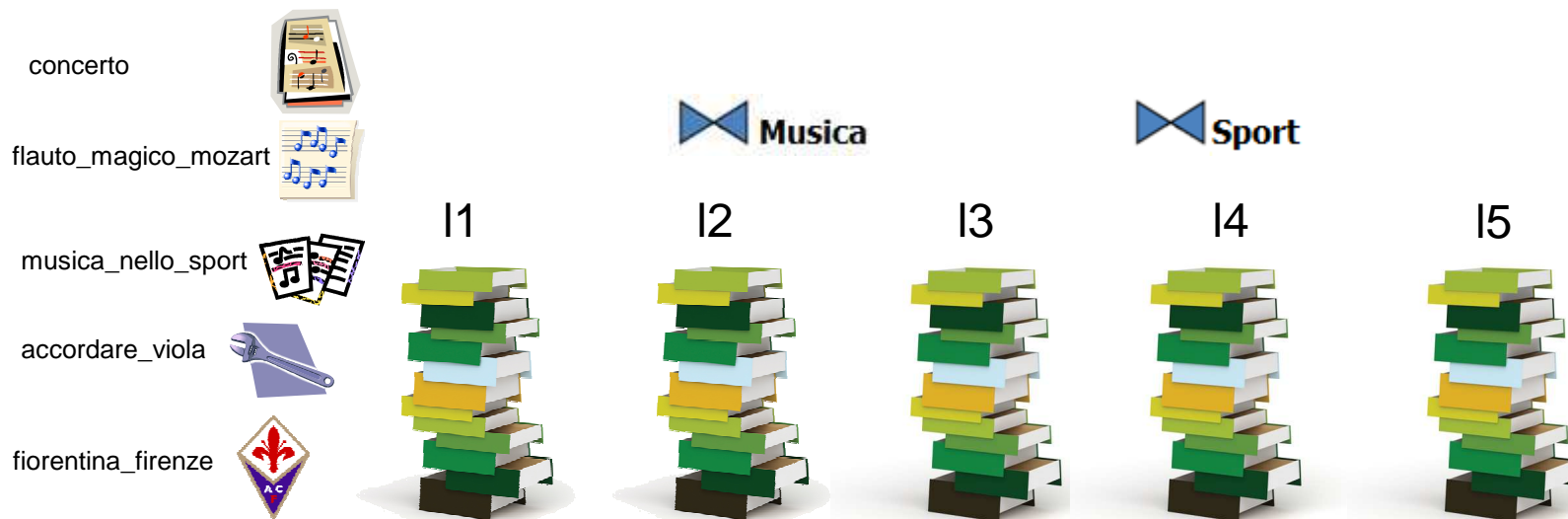


# Skip-and-Prune: Fase 1



9

- L'obiettivo è cercare di estrarre per ogni concetto i documenti più rilevanti nel contesto dell'utente, senza considerare ancora la query!
- Esecuzione degli operatori di Ranked Join su tanti stream di documenti ordinati quante sono le parole chiave rilevanti per il concetto a cui sono associati



# Skip-and-Prune: Fase 1

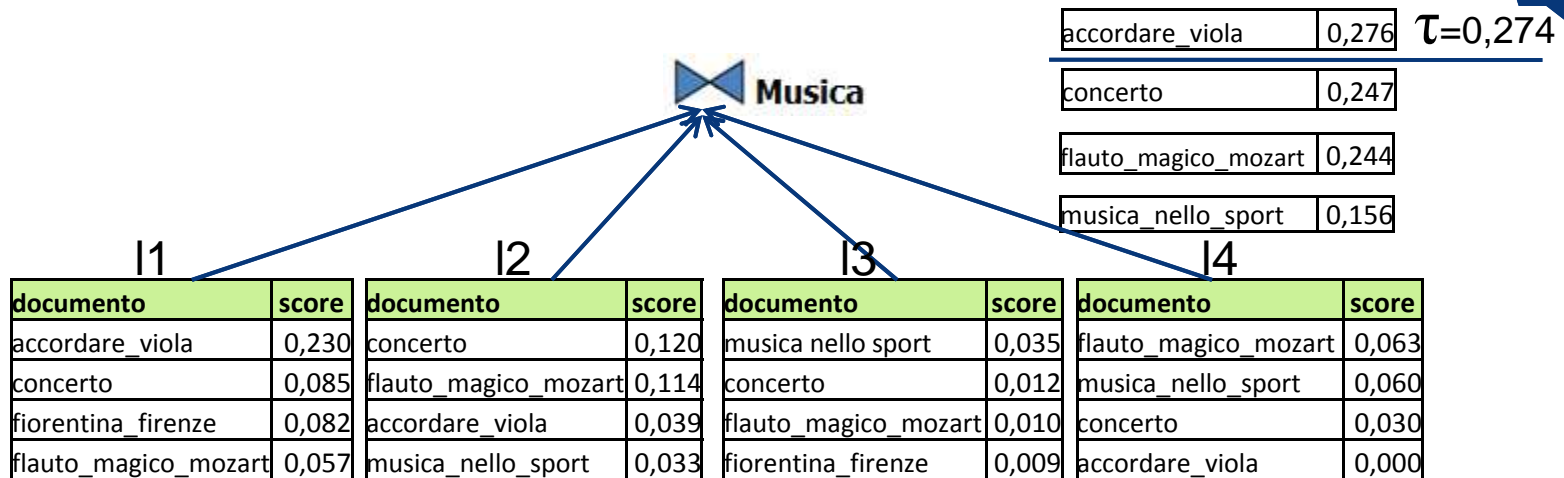


10

U	I1	I2	I3	I4	I5
musica	0,28	0,30	0,12	0,30	0,00
sport	0,10	0,00	0,20	0,00	0,70

D	I1	I2	I3	I4	I5
concerto	0,30	0,40	0,10	0,10	0,10
flauto_magico_mozart	0,20	0,38	0,09	0,21	0,00
musica_nello_sport	0,10	0,11	0,30	0,20	0,20
accordare_viola	0,81	0,13	0,06	0,00	0,00
fiorentina_firenze	0,29	0,00	0,08	0,00	0,40

Fase 2



# Skip-and-Prune: Fase 1



11

U	I1	I2	I3
musica	0,28	0,30	0,12
sport	0,10	0,00	0,20

## ATTENZIONE:

non abbiamo ancora considerato la query, stiamo solo inviando alla fase 2 i documenti più rilevanti di ogni concetto

## IDEA:

lo score di un documento è la sua coordinata nella dimensione di quel concetto



U	I1	I2	I3	I4
flauto_magico_mozart	0,244			
musica_nello_sport	0,156			
flauto_magico_mozart	0,063			
musica_nello_sport	0,060			
concerto	0,030			
accordare_viola	0,000			
concerto	0,012			
flauto_magico_mozart	0,010			
fiorentina_firenze	0,009			
musica_nello_sport	0,033			
concerto	0,120			
flauto_magico_mozart	0,114			
accordare_viola	0,039			
musica_nello_sport	0,033			
concerto	0,230			
flauto_magico_mozart	0,085			
accordare_viola	0,082			
musica_nello_sport	0,077			

# Skip-and-Prune: Fase 1



12

- L'operatore di Ranked Join restituisce i documenti che ottengono un punteggio oltre una certa soglia denominata  $\tau$
- La soglia tiene conto sia di D che di U consentendo di avanzare ai documenti con score di rilevanza per quel concetto maggiore di essa. Si aggiorna ad ogni Ranked Join
- I documenti in fondo alle InvertedLists sono meno importanti per quel concetto e, a meno che il loro ID sia all'interno di un insieme detto SkipSet che viene riempito dai documenti eliminati nella fase 2 (retroazione!), supereranno più tardi la soglia per accedere alla fase successiva

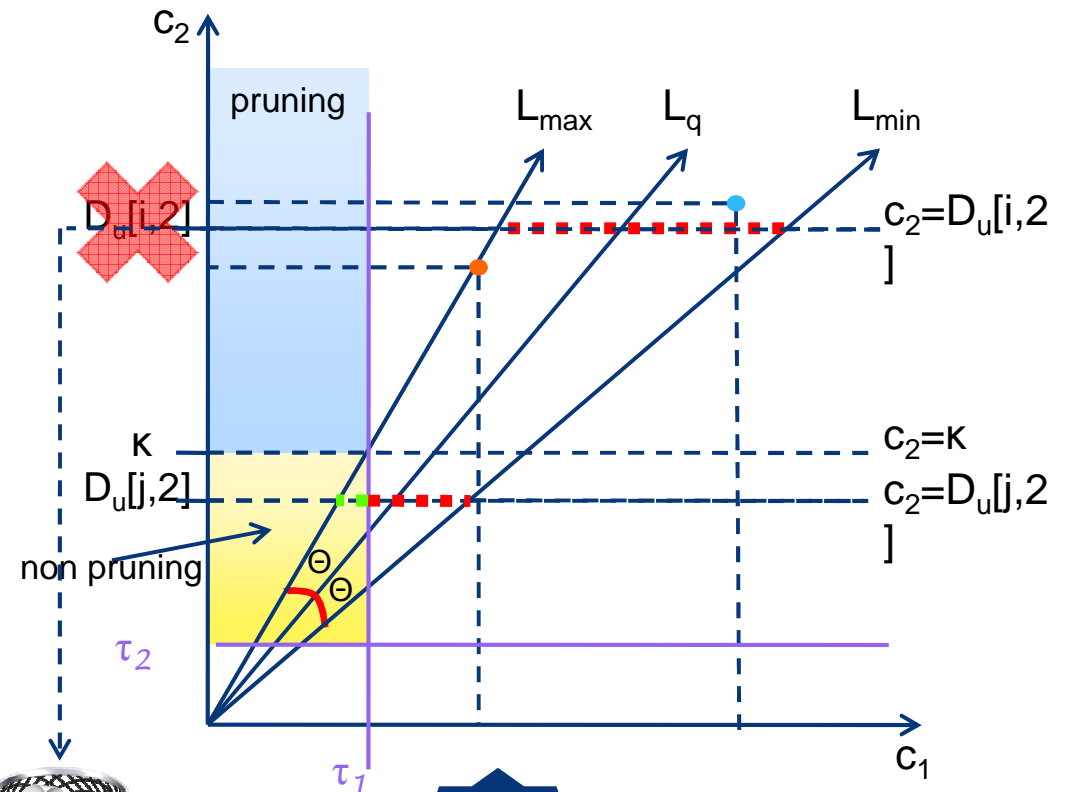


# Skip-and-Prune: Fase 2



13

- Spazio bidimensionale dei concetti
- La query dell'utente considerata nel suo specifico contesto è rappresentata con il vettore  $L_q$ . Obiettivo di questa fase è scegliere documenti vicini alla query nel piano
- Ogni documento uscente dalla fase 1 ha associato il punteggio maggiore alla soglia, rappresentante la coordinata nell'asse del concetto a cui è associato il Ranked Join. Se il documento ha superato entrambi i Ranked Join è un punto nel piano
- Supponiamo di considerare una query top-2; ricevuti i primi due documenti posso determinare i vettori  $L_{max}$  ed  $L_{min}$  che definiscono il triangolo dei risultati accettabili
- Consideriamo un nuovo documento  $d_i$  che supera il Ranked Join di  $c_2$ : ha una sola coordinata nello spazio e  $c_1$  è incognito.  $d_i$  è accettabile solamente se è posizionato sulla linea tratteggiata
- Sia  $K$  l'ordinata del punto dove  $L_{max}$  interseca  $\tau_1$ ;  $d_i$  può essere eliminato dato che la fase 1 limita le sue ascisse a valori inferiori a  $\tau_1$ : non potrà mai posizionarsi nella zona tratteggiata! Il suo id viene inserito nello SkipSet ed il Ranked Join di  $c_1$  non lo considererà



Skip Set



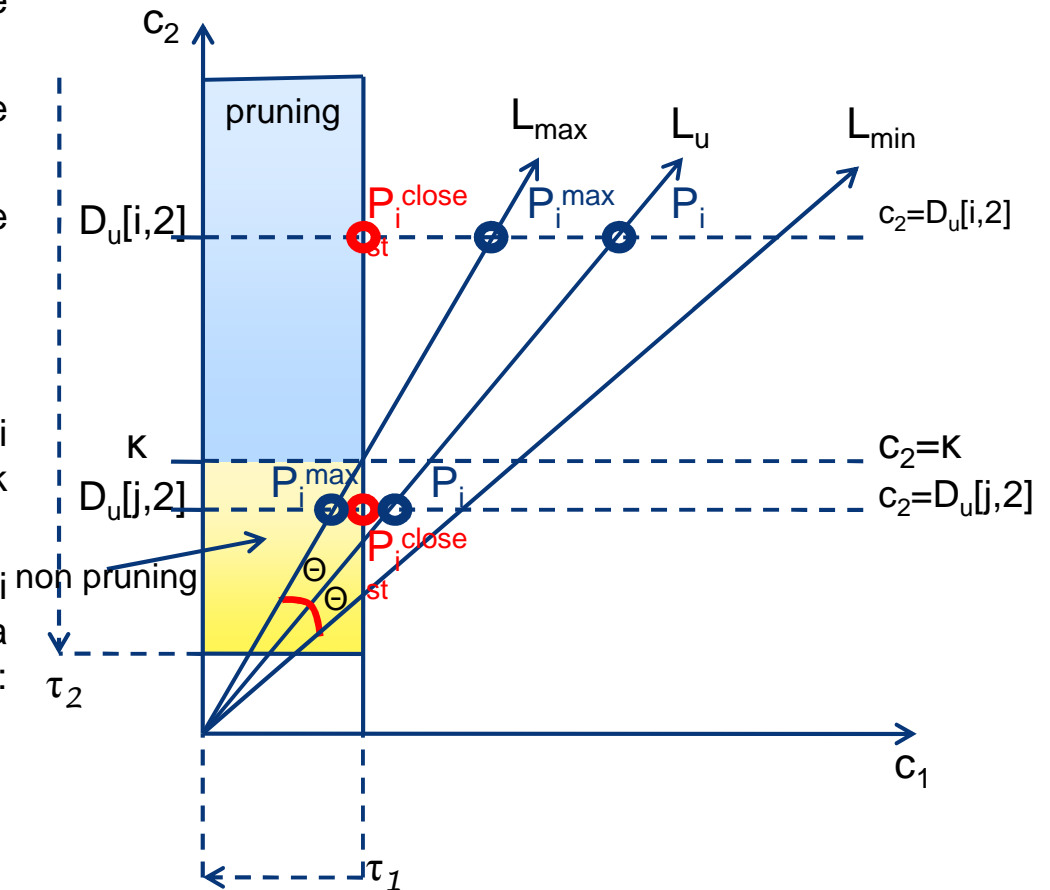
# Skip-and-Prune: Fase 2

## Criteri di Pruning



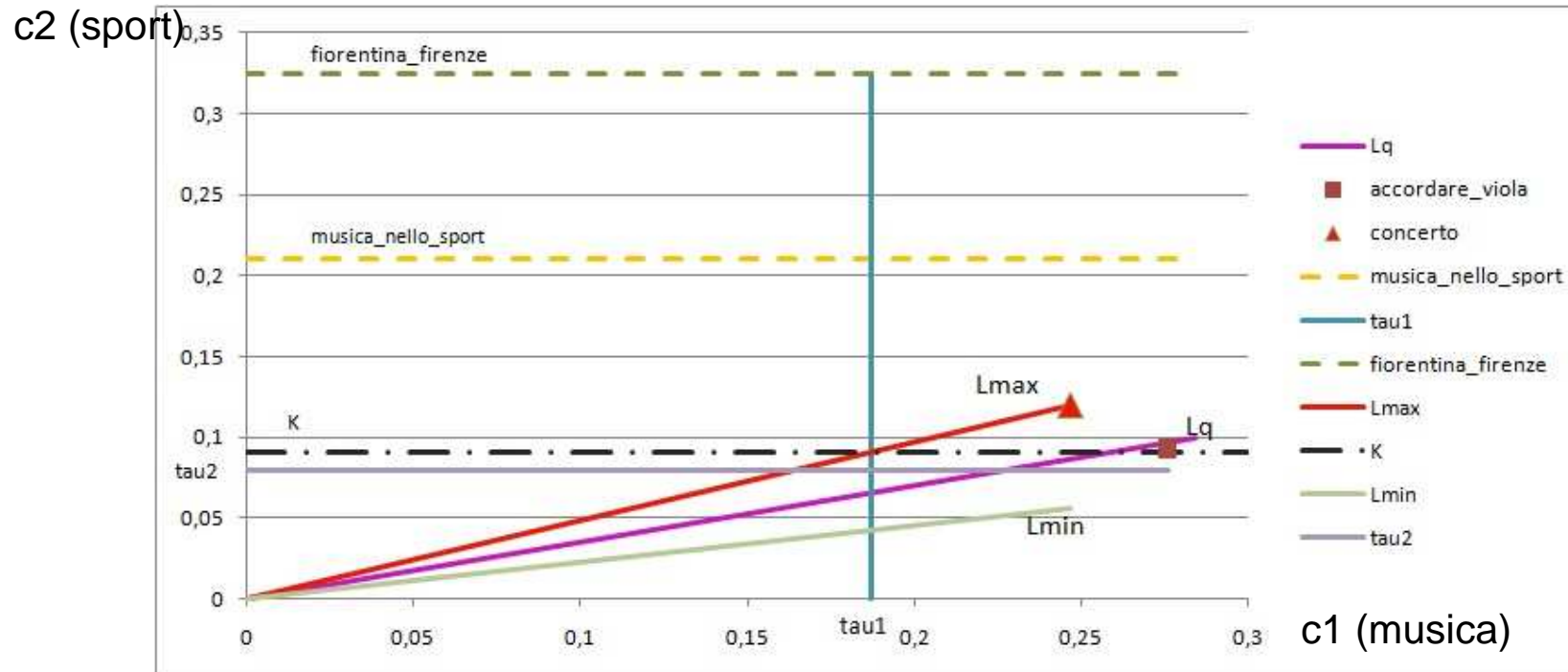
14

- $P_i$  rappresenta la query in corrispondenza della retta del documento  $H_i$ , tracciata grazie all'ordinata ottenuta dalla fase 1
- $P_i^{closest}$  rappresenta l'ascissa massima che  $D_u[i,2]$  potrà avere su  $H_i$
- $P_i^{max}$  è il punto limite accettabile per far sì che  $d_i$  sia tra i correnti top-k candidati
- $d_i$  verrà eliminato se:  
 $\Delta(P_i^{closest}, P_i) \geq \Delta(P_i^{max}, P_i)$
- Altrimenti il suo score sarà maggiore di  $min\_score$ , il più basso tra quelli dei top-k documenti
- $min\_score$  consente di fare le considerazioni geometriche: la regione accettabile è definita dal triangolo tramite l'equazione:  
 $\cos(\Theta) = min\_score$



# Che fine ha fatto l'esempio??

15



- Sui documenti “fiorentina\_firenze” e “musica\_nello\_sport” viene fatto pruning → il loro ID viene inserito nello SkipSet
- Come atteso, i documenti **migliori** sono “concerto” e “accordare\_viola”





# Skip-and-Prune: ulteriori raffinamenti...

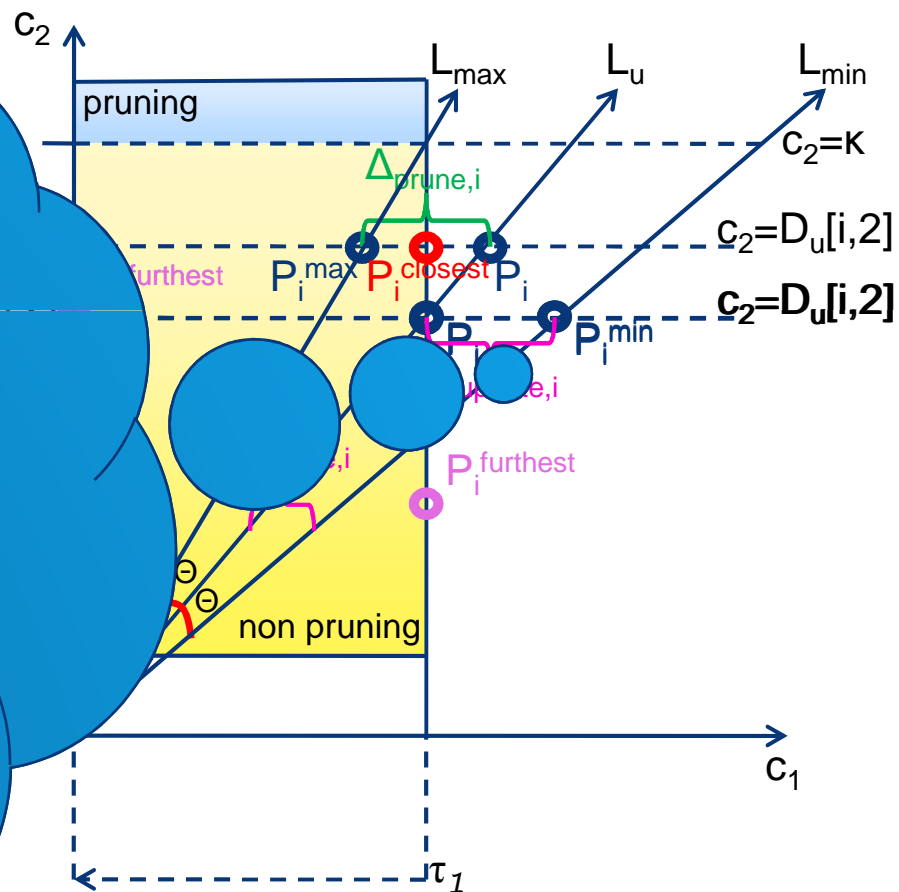
17

- Il raffinamento proposto consiste nel ridefinire  $min\_score$  di modo che possa riflettere i punteggi dei documenti visti solo parzialmente
- Prendiamo il documento  $d_i$  e i suoi  $c_1$  dei top-k col

L'efficienza dell'algoritmo di pruning dipende da quanto stretti sono i valori della distanza di pruning

$$\Delta_{prune,i} = \Delta(P_i^{max}, P_i)$$

per i documenti candidati. Questo valore dipende da  $min\_score$ , che, specialmente negli istanti iniziali del processo, sarà troppo basso per l'assenza di documenti completi



26/05/2010

# Skip-and-Prune: considerazioni

18

- Il meccanismo dello *skip set* assicura che i documenti vengano rimossi il prima possibile dagli stream di input per avere un minore tempo di elaborazione totale
- A differenza di TA, questo processo non può terminare finché gli InvertedFiles non siano stati completamente scanditi, a causa della non monotonicità della funzione coseno



# Valutazioni Sperimentali

19

- Sono stati fatti esperimenti per valutare l'**efficienza** dell'approccio proposto, confrontando i risultati con le soluzioni precedenti:



- ▣ **Scan-and-Choose (SC)**: scandisce tutti i documenti, li reinterpreta secondo il contesto dell'utente restituendo i migliori  $k$



- ▣ **Accumulator-based Inverted File (AIF)**: utilizza InvertedFiles e accumulatori per calcolare il punteggio globale di ogni documento



- ▣ **NaiveRanked Processing (NRP)**: simile a  $SnP$  perché considera solo i documenti rilevanti per ogni concetto, ma utilizzando una funzione di scoring monotona

# Valutazioni Sperimentali: dati sintetici

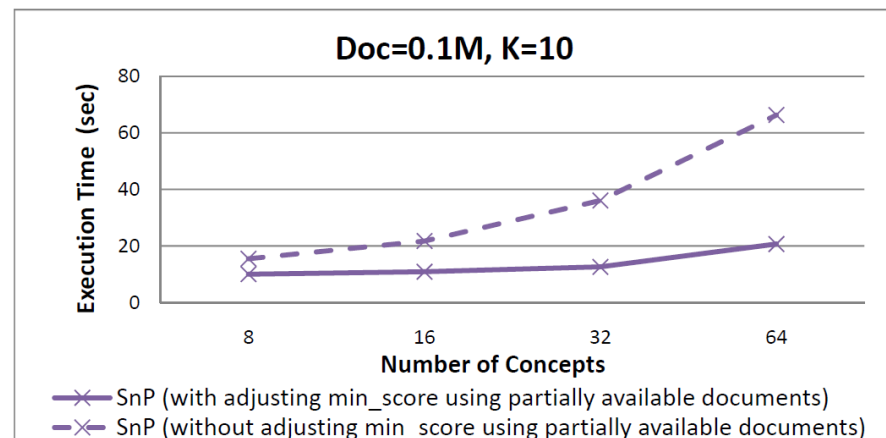
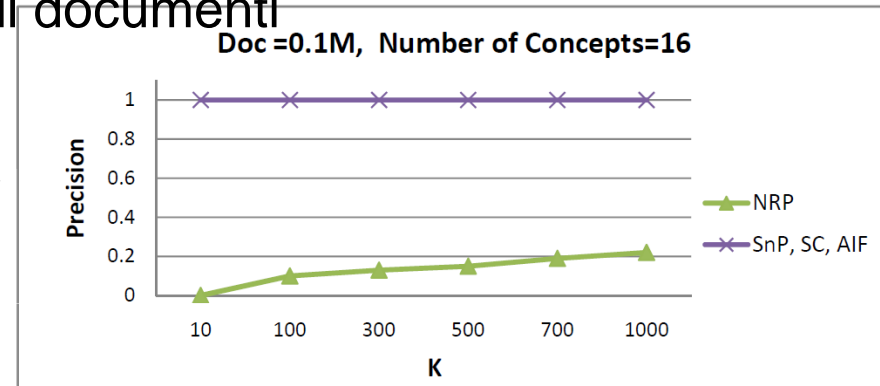
20

**Dati sintetici** generati sistematicamente, variando parametri come  $k$ , il numero dei concetti e il numero di documenti

Confronto della precisione media dei top- $k$  risultati in funzione di  $k$ , numero dei documenti restituiti all'utente.

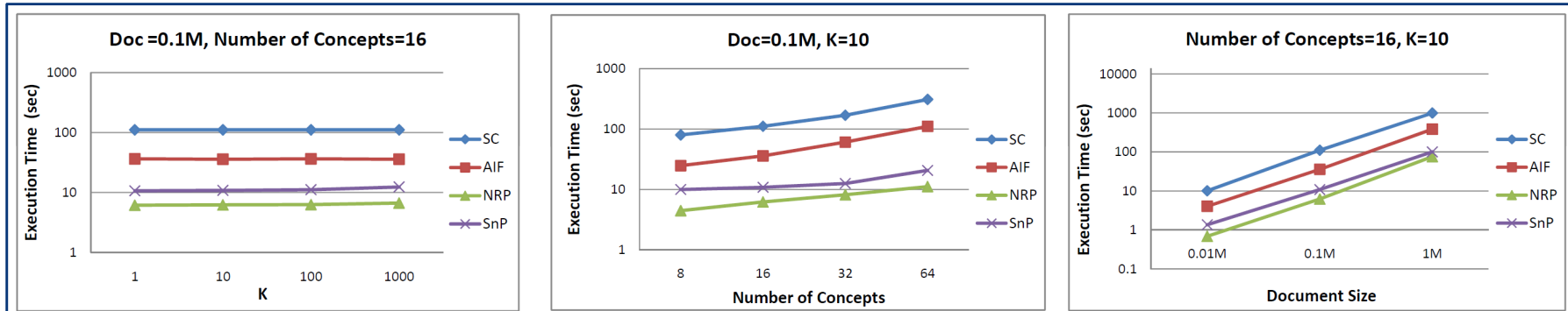


Impatto dell'aggiustamento del  $min\_score$  considerando documenti visti parzialmente

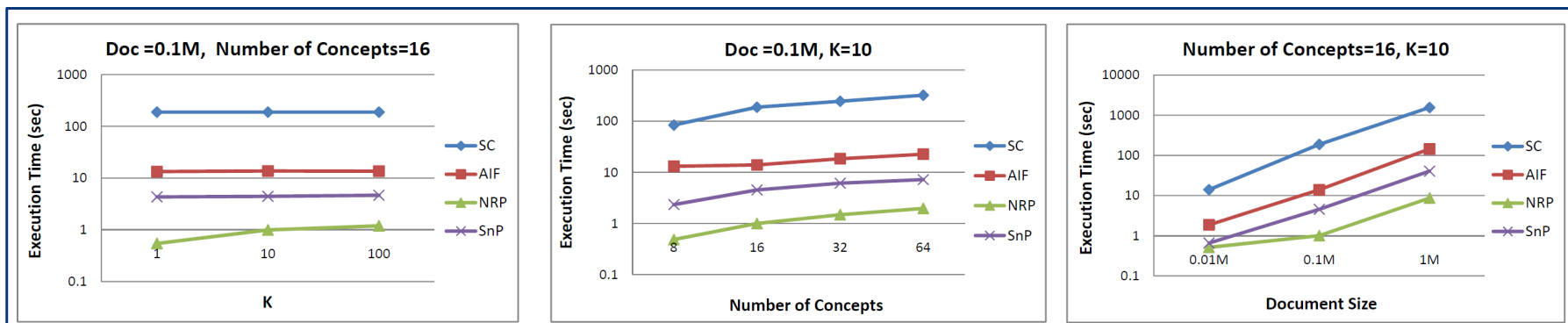


# Valutazioni Sperimentali: dati sintetici

21



Matrici D dense: tempi di esecuzione al variare di differenti parametri



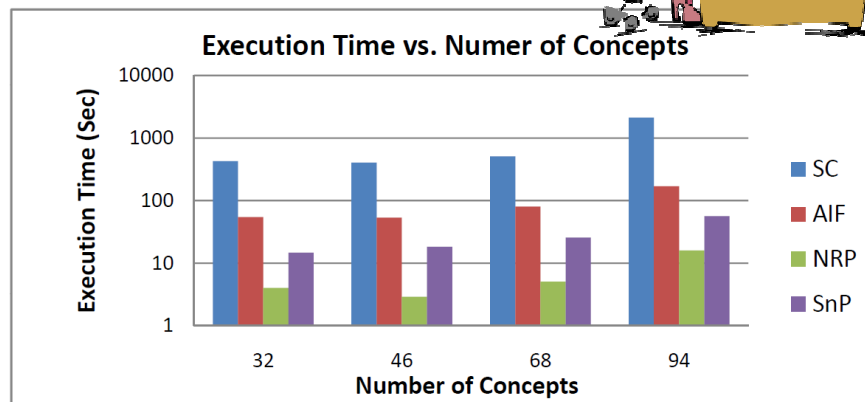
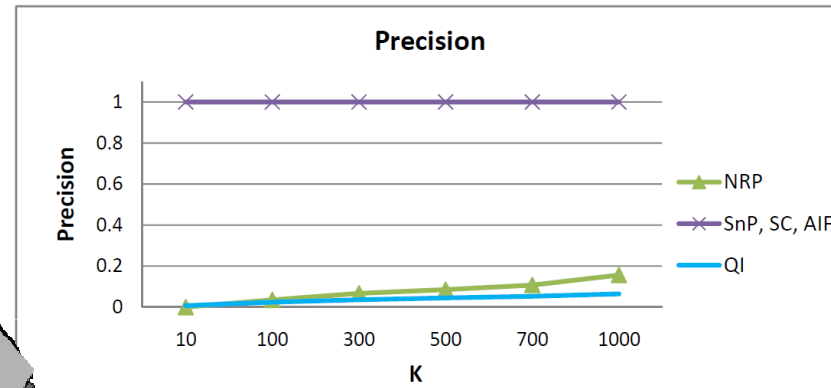
Matrici D sparse: tempi di esecuzione al variare di differenti parametri

# Valutazioni Sperimentali: dati reali

22

**Dati reali** estratti da ACM digitallibrary e ScienceDirect

Precisioni medie dei top-k risultati sui dati di insiemi reali. Ulteriore algoritmo di confronto: **Query-only Interpretation**



Gruppo 14

Tempi di esecuzione dei differenti algoritmi al variare del numero di concetti significativi.

L'algoritmo SnP qui proposte di uno o due ordini più veloce di SC e AIF

26/05/2010

# Conclusioni

23

- L'algoritmo qui presentato, consente efficientemente di calcolare query Top-k in presenza di un contesto utente con una funzione di scoring non monotona basata sul coseno



## Vantaggi principali:

- ➔ È efficiente e scala molto bene al crescere della dimensione dei documenti e del numero di dimensioni nello spazio dei concetti
- ➔ Ha tempi di esecuzione migliori degli algoritmi precedenti dato che consente di eliminare documenti non inerenti al contesto e alla query
- ➔ Fornisce risultati ad alta precisione

- Estensioni future possibili ad altre funzioni di scoring non monotone

24

FINE

**GRAZIE !**

*Votagruppo 14 !!*