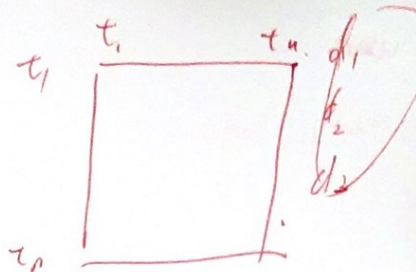


$$\vec{q} = \begin{matrix} t_1 & t_2 & t_3 \\ 0.7 & 0.6 & 0.1 \end{matrix}$$

$$\vec{q}_1 = \begin{matrix} t_1 & t_2 & t_3 & \dots \end{matrix} \quad \text{correlates.}$$



can have some subtopic \rightarrow can suggest users.

\downarrow
disambiguate query.

if $\vec{q} = t_1$. java = C++
 programming }

coffee
espresso. }

info. several

① correlated t . select the best correlated terms.

② diversity t . \rightarrow help user to make query more specific.
 \downarrow
try to find an the trade off between this two.

If I choose k top k correlated terms, how could I
say the diversity.

users

query log.

user id

1 $t_1 t_2 t_3$

2 $t_1 t_2$

1 $t_2 t_3 t_4$

2 $t_4 \longrightarrow$ a related topic.

:

x $t_3 t_4 \dots$

Clustering IR.

If a query \rightarrow cluster

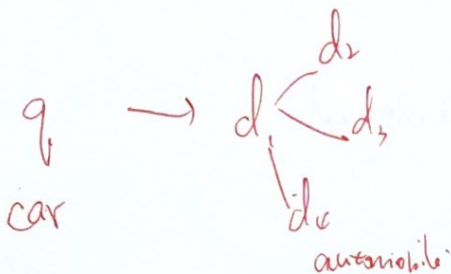
term-term matrix

\downarrow

query-cluster matrix.

Search Result clustering \rightarrow increase the diversity

Scatter-clustering \rightarrow visualise.



Flat algorithm



k-means.

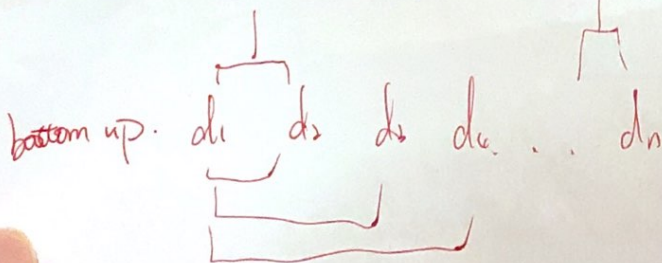


Hierarchical algorithm

Collection

sports

~~politics~~ politics



IR-3.

N docs and K clusters.

~~K~~ ^{K} -means. \rightarrow use human judgement for evaluation.

K -NN: \rightarrow classification, use prior knowledge in algorithm.

size of each cluster should be balanced.

①. RSS. (interior distance)

②. K . (make ~~it~~ a tight cluster)

③. give penalty to the cluster which size is too small or too large.

~~$dis = \sum dis_i - k^2$~~