# Distributed Index Generation (MAP-REDUCE)

## Phase 1: Build the Term-Partition Index
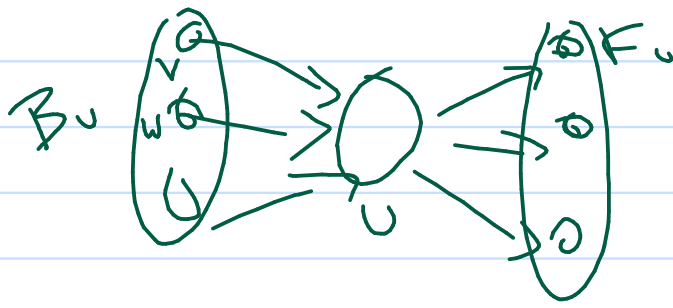
Input: Documents

Master

**Map**

**Parsers**

P1
P2
P3

**Reduce**

Inverters

| A | G | Q |
| A | G | R |
| A | G | R |

Inverter
Inverter
Inverter

A
G
Q

Inverted Index with postings    Postings

## Phase 2: Build Document Partition Index

Input: Term-Partitioned Postings (Output from Phase 1)

**Parsers**

P1
P2
P3

**Inverters**

I 1
I 2
I 3

0 - 10k

10,001 - 20k

20,001 - 30k

Document Segments

With Doc Partitioning,
the entire query is sent to every node which reduces (alleviates)
the doc list merging that takes place with Term-Partition indexes.

## Nov 1

## EXAM Result

## Use IDF for Query Calculation

Nov 5 — Talked to Dr. D., the question indicated which fmla to use
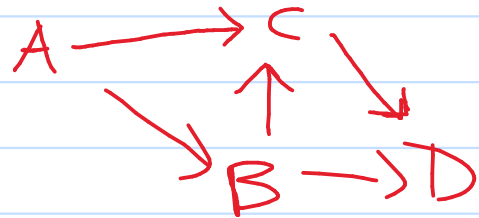"do not normalize the query" means to leave it out of the Cosine Similarity fmla altogether.

$$\frac{R(v)}{4} + \frac{R(u)}{2}$$

## Build an Adjacency Matrix

①

$$\begin{pmatrix} & A & B & C & D \\ A & 0 & 1 & 1 & 0 \\ B & 0 & 0 & 1 & 1 \\ C & 0 & 0 & 0 & 1 \\ D & 0 & 0 & 0 & 0 \end{pmatrix}$$



P

$$\begin{pmatrix} & & \diagdown & \\ 0.25 & 0.25 & 0.25 & 0.25 \end{pmatrix}$$

Since D has no OUT links, just assign the Prob. that they jump to any node.

$$P_r = \frac{1}{N} = \frac{1}{4} = 0.25$$

Then fill in the other rows based on link probabilities

$$\begin{pmatrix} 0 & .5 & .5 & 0 \\ 0 & 0 & .5 & .5 \\ 0 & 0 & 0 & 1 \\ .25 & .25 & .25 & .25 \end{pmatrix}$$

Prob of user following a link

then... $P = (1-\alpha)P$ (for all but last row)

$$+$$

②   $(\alpha)(\frac{1}{N})$   (except last row b/c already calculated)

③ P.R.

$\vec{X}_0 = (1,0,0,0) \leftarrow$ can choose any starting vector.
$\vec{X}_1 = (\vec{X}_0)(P)$
$\quad = (.05, 0.45, .45, .05)$
$\vec{X}_2 = \vec{X}_1 \cdot P$
⋮

Until the $\Delta$ is very small. -- then we have (P.R.)

─────────────────────────

# Nov 15

①

## Adj Matrix

$A \rightarrow B$
$\quad \updownarrow$
$\searrow C$

$\begin{array}{c} A \\ B \\ C \end{array} \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$

Ⓐ Page Rank

- all have outlinks!, unlike last class's example

$P_r = 1/\text{sum}(1's)$

$\begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \leftarrow$ Pr of user following a link

$\alpha$, say, $0.1$
So P: $0.9$ $\Rightarrow$ $\begin{pmatrix} 0 & .45 & .45 \\ 0 & 0 & .9 \\ 0 & .9 & 0 \end{pmatrix}$

# Cont'd...

Pr of teleport: $\to \frac{1}{3}$

$$Pr_{(Link)} + (1-\alpha)\frac{1}{3}$$

$$= \begin{pmatrix} .03 & .48 & .48 \\ .03 & .03 & .93 \\ .03 & .93 & .03 \end{pmatrix}$$

Finally, choose $\vec{X_0}$

$$\vec{X_0} = (1, 0, 0)$$

$$\vec{X_1} = \vec{X_0} \cdot P$$
$$\vec{X_2} = \vec{X_1} \cdot P$$
$$\vdots$$

Until convergence
$$\vdots$$

or after so many steps.

---

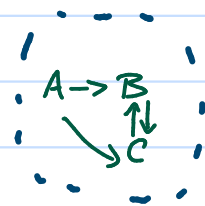# Hits Algorithm

Need Hub score + auth score!

1. Intial vectors:
$$\vec{h_0} = (1, 1, 1)$$
$$\vec{a_0} = (1, 1, 1)$$

A → B
↑↓
↘C

2. Check in-/links. —— B & C have 2 in-Links

A has none ∴ ↗

$$\vec{a_1} = (\overset{\curvearrowleft}{0}, (1+1), (1+1))$$
$$= (0, 2, 2)$$

Normalize :→ *Make so that the MAX value is 1*

$$a_1 = (0, 1, 1) \qquad [0, \frac{2}{2}, \frac{2}{2}].$$
divide by largest

Hub Score: Sum of Auth. Score:

$$h_1 = (2, 1, 1)$$
Sum of outbound nodes' auth scores

$$\underset{NORmal}{h_1} = (1, .5, .5)$$

$$a_2 = (0, 1, 1) - \text{converged,}$$
So stop.

✳ authority score subsequently is Sum of inbound hub scores.

# Recommender Sys Cont'd

### Movies

|      | M1 | M2 | M3 | ... |
|------|----|----|----|-----|
| Alice | 3 | 2 | 5 |     |
| Bob  | 4  | 5  | 2  |     |
| C    | ∅  | 3  | 4  |     |
| ⋮    |    |    |    |     |

Recommend movies to users } diff results!
or users to movies

empty set → hasn't seen

## Content-based

Utility $u(c,s)$;

Similar

$u(c, s_i) \mid s_i \in S$

### User
$\overrightarrow{W_c}$ — of $W_{ck}$ ...keywords -- Need to determine these!

avg rating of...

### Content
$\overrightarrow{W_s}$

directors
animator
drama

**Need to do Normalization**

Bob (Profile) : $(.9, .4, .2, 0, .3, .9)$

Predict Rating for "UP"

$MDD: (4+5)/2 = 4.5 \rightarrow$ Normliz $\rightarrow 0.9$

$2/5 \rightarrow 0.4 \quad \}[1]$

$1/5 \rightarrow 0.2 \quad \}[2]$

Up : $(0,0,0,1,1,0)$

Sim(Bob, Up)

$= 0.3 \times \underline{1} = 0.3$
$\quad\quad\quad (1.8)$

DarkKnight = $(1,0,0,0,0,1)$

# User-based CF

Alice avg $(3,2,5,4) = 3.5$ } only commonly
Bob avg $(4\ 5\ 2\ 1) = 3$    }        rated

$$\overset{m1}{(3-3.5)}\overset{m2}{(4-3)} + (2-3.5)\overset{m3}{(5-3)} + \cdots \qquad = -6$$

↑        ↗
A rating   A's avg

$$\frac{-6}{\sqrt{\sum component^2 s}} \quad \to always\ -1\ to\ 1$$

?

∴ $= -0.85$

highest is
Best!

$$Rating_{BOB\ UP} = 1/\left(\sum \underset{People}{(high\ similarity)}\right) * \left(\sum (SimPerson)(SimPerson\ Rating)\right)$$

# Item-based
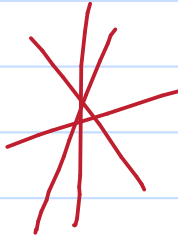
Clustering   2D Clusters →

Choose a
Prototype centroid

## K-mean

$$\vec{\mu}(c) = \frac{1}{|c|} \sum_{\vec{x} \in c} \vec{x}$$

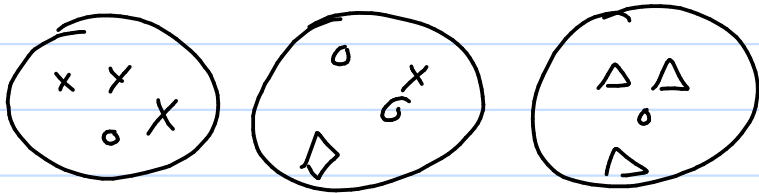- Select $h$ docs randomly as Centroids

# Hierarchical

- at the bottom level, every object in its own cluster; same at the top.

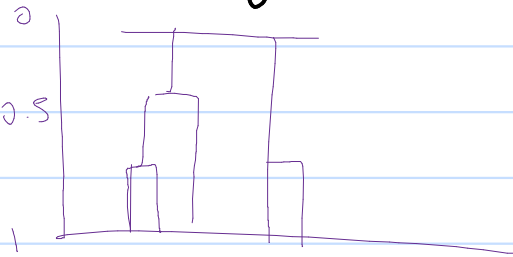Initial centroid selection affects the resulting clusters.

## Purity EX

$$Purity = \frac{\sum (majority\ items\ in\ each\ Circle)}{N}$$

$$Rand\ Index = \frac{TP + TN}{TP + FP + FN + TN}$$

# HAC

- Repeatedly joins two clusters until there is only 1.

- A Dendogram



## Sim



complete link

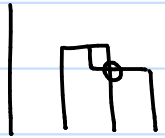closest pair ("single link")

3 - centroid. (4)

4. Avg (6 links)

## Single Link Clustering — find closest cluster by measuring from inner side of clusters.

- dist betw. clusters is closest pairs
- group closest things repeatedly
- results in undesirable long chains

# Complete Link   — measure from outer edge
 - Longest dist   — better balance; instead of having lots of clusters
                    with only 1 item
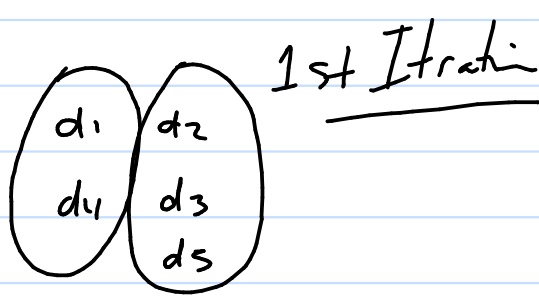                  — downside: outliers screw up the grouping

Inversion in Dendogram    ← don't use controid
                                                 HAC for this reason

---

Q3 HW

|     | T1 | T2 | T3 |
|-----|----|----|----|
| D1  | :  | :  | :  |
| :   |    |    |    |
| D5  | .  | .  | .  |



K-Means first

1) k=2;  controids d1, d2
         decide the rest; which cluster....

$d_3: d_{1,3} = [inner product] 1\cdot 0 + 0\cdot 9 + 0\cdot 4 = \emptyset$
$\quad d_{2,3} = 0.8$  (d2) ✓

$d_4: d_{1,4} = .9$ ✓  (d1)
$\quad d_{2,4} = .3$

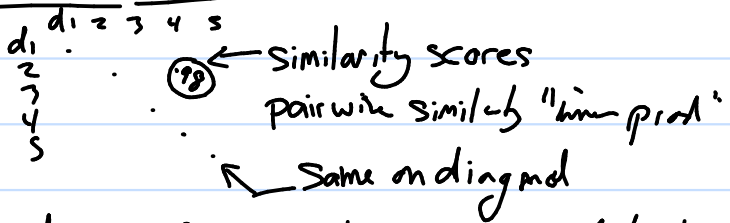$d_5: d_{1,5} = 0.6;  d_{2,5} = 0.98$  (d2)



1st Iteration

2) Now calc new centroids
$C_1 = avg(d1, d4) = (0, 0.95, 0.2)$
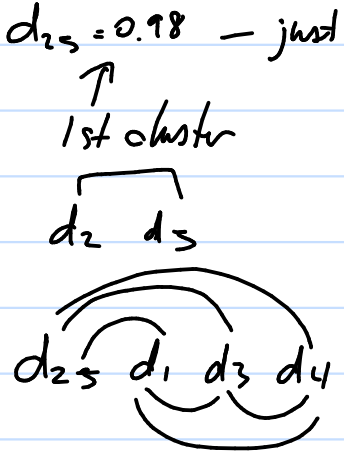$C_2 = \ldots\ldots \quad (0.83, 0.23, 0.37)$

Calc distances again:
$d_{1_{C_1}} = ((.9)(.95) + 0.00) = .935$  ?
$d_{1_{C_2}} = (.9)(.23) + (.4)(.3)) = 0.3$ ? (smaller)
do calcs for the rest of the docs to put in either $C_1$ or $C_2$

don't include
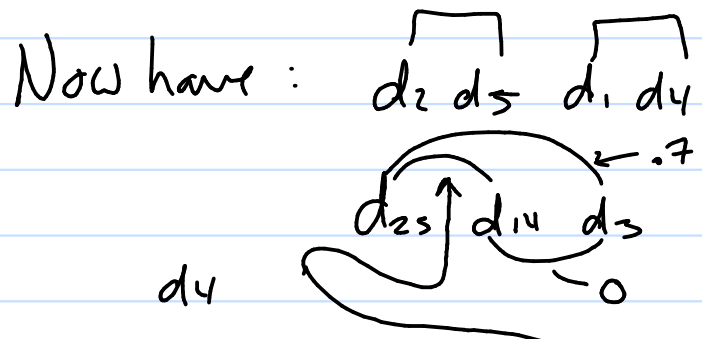the controid

in subsequent iterations.

# HAC EX!

$d_i$
|     | $d_1$ | 2 | 3 | 4 | 5 |
|-----|-------|---|---|---|---|
| $d_1$ |  | · |  | (.98) ← Similarity scores |  |
| 2 |  |  | · |  |  |
| 3 |  |  |  | · |  |
| 4 |  |  |  |  | · |
| 5 |  |  |  |  |  |

← Similarity scores
Pairwise similarity "inner prod"

↖ Same on diagonal

$d_{23} = 0.98$ — just do 1 corner of Matrix

↑
1st cluster

$d_2 \quad d_3$

$d_{23} \quad d_1 \quad d_3 \quad d_4$

**Find → after Midterm, focus on Web Search!**

$Sim(d_1, d_{23}) = ? \rightarrow$ use Complete Link
→ farthest
→ Least Similar!
→ Smaller → 0.47

$d_3, d_{23} = 0.7$
$d_4, d_{23} = 0.3$

Then find most similar to cluster → $d_1 \& d_4$ (b/c 0.9)

Now have: $d_2 \, d_3 \quad d_1 \, d_4$

$d_{23} \quad d_{14} \quad d_3$ ← .7

→ 0

|       | $d_1$ | $d_4$ |
|-------|-------|-------|
| $d_2$ | .47   | .3    |
| $d_3$ | .6    | .4    |

} pick lowest → 0.3

(Smallest of 6) from Matrix

0.3

.7
.98
.9

$d_2 \quad d_3 \quad d_3 \quad d_1 \, d_4$

Can calc top-most score too → 0.3

For Single Link; still start at $d_{23}$ (0.98), but now choose highest scores.
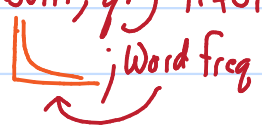
# Exam

"Item-Based" CF; same as assignment (the values anyway)
- avg of items rather than users.
- basically the vector will be in columns instead of rows
- Finds most similar ITEMs instead of people.

- Know K-means + hierarchical

# Review

- Lectures after MID.
- UI -> qry in; qry results; qry reformulation
- Docs & Queries -> Zipf ⌐_ ; Word freq.
- Compression:
        - Decompression speed is most important!
Query Intent (Broder): navigational, informational, transactional

---

(NB) Web Search ("Focus on this")                    (NB)
  1. W.search Basics
     Key differences: lints, query context, users/docs, spaming, advertisements
  · how to estimate Index size
  ·  "  "  detect near duplicates
  · ranking signals -> ① Content ② Link (PR) ③ Usage [clicks]
2. Crawler
   - must have features: robust, politeness | Should-have: efficient, etc
   - Crawl process: seed set -> fetch -> parse -> extract links & text, dupl. check -> URL frontier
   - Scheduling
     Architecture ("Mercator")
3. Link Analysis => PR, Hits

# Recommender

- Content-based, collab. filtering
  - memorize fmla's
  - Know limitations of each
  -

# Clustering
                    ② Recalc centroids

K-means: ① Initial seeds ③ Iterate till objective function is optimized

Hierarchical Aggl. Clust. (HAC)
  - group repeatedly till only 1 cluster.
    - Single link (measure from near side)
    - Complete (measure from far side)
  - Evaluation - ext criterion → purity, rand. index.