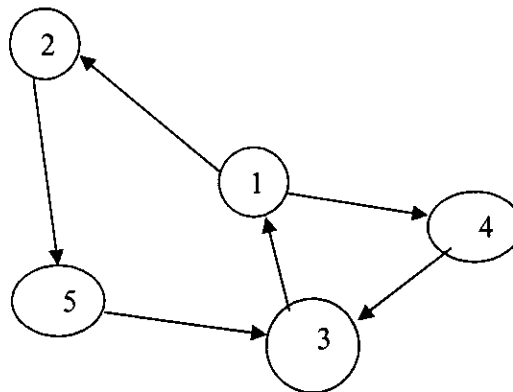


אוניברסיטת בן-גוריון - המחלקה להנדסת מערכות מידע
קורס איחזור מידע וספריות דיגיטליות
סמסטר אביב תשס"ו - 09.07.06 – מועד א

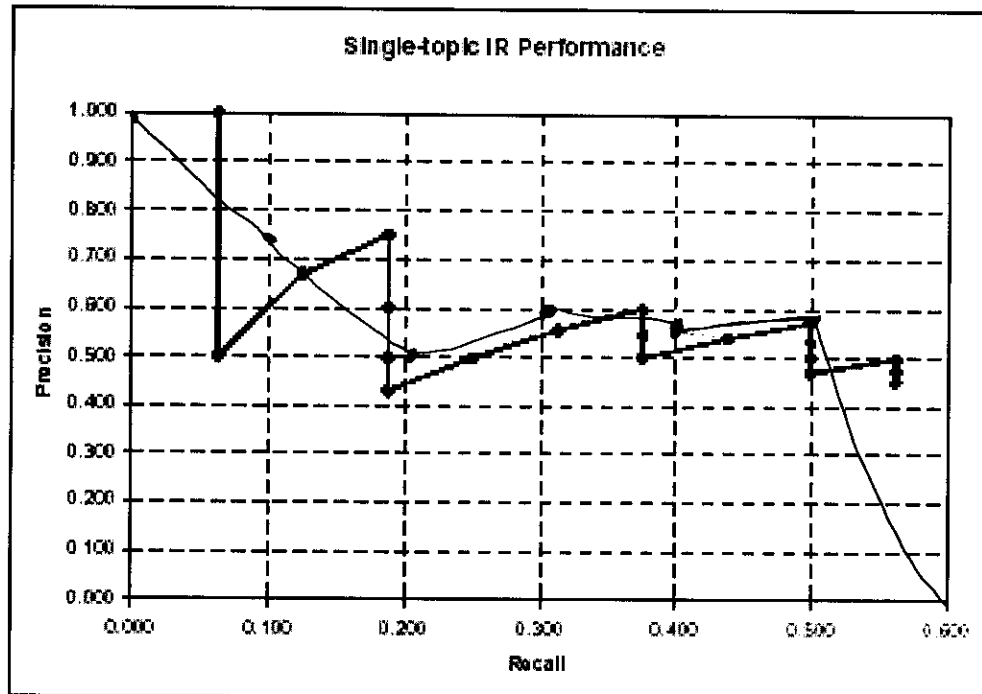
מרצה: ד"ר ברכה שפירא, אסיסטנט : ארז שלום
 משך המבחן: 3 שעות, חומר עזר מודפס או כתוב - מותר, מחשבון- מותר
 יש לענות על כל 5 השאלות. יש להתזיר את השאלונים.

1. (15%) נתונה רשת (בציור) המתארת צמתי Web והצבעות ביניהם :



- א. (3%) אם נוסף קישור מצומת 2 לצומת 4, האם וכיצד יושפעו ערכי HUB ו- AUTHORITY (על פי אלגוריתם HITS) של צומת 5 (כלומר, האם הערכים יהיו גבוהים, או נמוכים יותר, או ללא שינוי ומדוע).
- ב. (7%) חשב שתי איטרציות של pagerank של הגרף שבציור לאחר ההוספה מסעיף א. (ללא נרמול). הערך איזה צומת תהיה בעלת ערך pagerank גבוה ביותר לאחר ההתכנסות. הסבר את הערכתך.
- ג. (5%) האם אלגוריתמים של link-analysis כדוגמת pagerank מודדים גם את פופולאריות הדף אצל הגולשים. אם כן, כיצד? ואם לא איך אפשר למדוד פופולאריות של דפים אצל גולשים?

2. (30%) הגרף הבא מתאר תוצאה של שאילתא אחת במנוע חיפוש שבה הוחזרו 20 תוצאות. ידוע שלשאילתא יש 16 תוצאות רלוונטיות (שלא כולן הוחזרו על ידי המנוע). הנקודות על הגרף מראות את ה-precision ו-ה-recall בכל נקודות המסמכים שהוחזרו. (כלומר כל נקודה מייצגת precision ו-recall למסמכים שהוחזרו עד נקודה זו על ידי המנוע).



- א. (6%) יש לצייר על גבי הגרף הקיים את הגרף בנקודות ה recall הסטנדרטיות באמצעות אינטרפולציה (בשיטת האינטרפולציה שנלמדה בקורס)
- ב. (6%) אילו מסמכים מבין המסמכים שהמנוע החזיר היו רלוונטים? – התשובה צריכה לכלול מספרי מסמכים לפי סדר דירוגם על ידי המנוע. (למשל אם התשובה היא -1,3,5 - הכוונה היא שהמסמכים שדורגו במקומות האלו על ידי המנוע הם רלוונטים).
- ג. (6%) חשבי את ה precision הממוצע (Mean Average Precision) (הראה את דרך החישוב)
- ד. (3%) חשבי את R-Precision
- ה. (3%) מהו ה Precision בנקודת 10 מסמכים?
- ו. (6%) תאר שתי חולשות של מאגרי הבדיקה המסופקים על ידי TREC.

(20%) בקובץ Posting של אינדקס הופכי מסויים כל כניסה כוללת: [docid, frequency of term]
נתון האינדקס של מאגר במבנה שהוזכר למעלה ובו 5 מסמכים:

Dictionary	Posting
Mother--→	{1,3},{3,2},{4,1},{5,3}
Father-- →	{2,4},{3,1},{5,2}
Family---→	{2,1},{4,3}
Brother--→	{1,4},{3,3},{4,2}

- א. (4%) אילו מסמכים יוחזרו ובאיזה סדר לשאילתות הבאות, במנוע חיפוש במודל בוליאני טהור:
Mother and Family (1)
Mother or Family (2)

5. In Authority "The Elf King" 5. In NUB - 10

1 μm	2 μm	3 μm	4 μm	5 μm
$\frac{1 \mu\text{m}}{0.15 + 0.85 = 1}$	$\frac{2 \mu\text{m}}{0.15 + 0.85/2 = 0.575}$	$\frac{3 \mu\text{m}}{0.15 + 0.85 + 0.85 = 1.85}$	$\frac{4 \mu\text{m}}{0.15 + 0.85/2 + 0.85/2 = 1}$	$\frac{5 \mu\text{m}}{0.15 + 0.85/2 = 0.575}$
$0.15 + 1.85 * 0.25 = 1.7225$	$0.15 + 0.85/2 = 0.575$	$0.85 + 0.85 * 1.575 + 0.15 = 1.48875$	$0.85/2 + 0.85/2 + 0.575 + 0.15 = 0.819$	$\frac{9 \mu\text{m}}{0.85/2 + 0.575 + 0.15 = 0.394}$

~~אם חלל נכנס לאזור אחר לא יושב: קיבלה את הקרן של 0.2 ש"ח~~

[illegible]

חברה
 1
 2
 3
 4
 5
 6
 7
 8
 9
 10
 11
 12
 13
 14
 15
 16
 17
 18
 19
 20
 21
 22
 23
 24
 25
 26
 27
 28
 29
 30
 31
 32
 33
 34
 35
 36
 37
 38
 39
 40
 41
 42
 43
 44
 45
 46
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65
 66
 67
 68
 69
 70
 71
 72
 73
 74
 75
 76
 77
 78
 79
 80
 81
 82
 83
 84
 85
 86
 87
 88
 89
 90
 91
 92
 93
 94
 95
 96
 97
 98
 99
 100
 101
 102
 103
 104
 105
 106
 107
 108
 109
 110
 111
 112
 113
 114
 115
 116
 117
 118
 119
 120
 121
 122
 123
 124
 125
 126
 127
 128
 129
 130
 131
 132
 133
 134
 135
 136
 137
 138
 139
 140
 141
 142
 143
 144
 145
 146
 147
 148
 149
 150
 151
 152
 153
 154
 155
 156
 157
 158
 159
 160
 161
 162
 163
 164
 165
 166
 167
 168
 169
 170
 171
 172
 173
 174
 175
 176
 177
 178
 179
 180
 181
 182
 183
 184
 185
 186
 187
 188
 189
 190
 191
 192
 193
 194
 195
 196
 197
 198
 199
 200
 201
 202
 203
 204
 205
 206
 207
 208
 209
 210
 211
 212
 213
 214
 215
 216
 217
 218
 219
 220
 221
 222
 223
 224
 225
 226
 227
 228
 229
 230
 231
 232
 233
 234
 235
 236
 237
 238
 239
 240
 241
 242
 243
 244
 245
 246
 247
 248
 249
 250
 251
 252
 253
 254
 255
 256
 257
 258
 259
 260
 261
 262
 263
 264
 265
 266
 267
 268
 269
 270
 271
 272
 273
 274
 275
 276
 277
 278
 279
 280
 281
 282
 283
 284
 285
 286
 287
 288
 289
 290
 291
 292
 293
 294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316
 317
 318
 319
 320
 321
 322
 323
 324
 325
 326
 327
 328
 329
 330
 331
 332
 333
 334
 335
 336
 337
 338
 339
 340
 341
 342
 343
 344
 345
 346
 347
 348
 349
 350
 351
 352
 353
 354
 355
 356
 357
 358
 359
 360
 361
 362
 363
 364
 365
 366
 367
 368
 369
 370
 371
 372
 373
 374
 375
 376
 377
 378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436
 437
 438
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468
 469
 470
 471
 472
 473
 474
 475
 476
 477
 478
 479
 480
 481
 482
 483
 484
 485
 486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519
 520
 521
 522
 523
 524
 5

צומח 3 מחבורה אדומה 1 ויקר עם את ה PR + 8% זאת היעדרות
הם בעיקר LN/L גלוקה.

1) link analysis היא מודדים באופן יחסי את הבוטומים של הדף
 אדם אחד. אדם אחד יחסיבא שבו בין הדפים באינטרנט. אדם אחד באופן
 שבו אדם אחד קושר ~~באופן~~ שבו באופן הדף שבו אדם אחד
 אדם אחד באופן הדף.
 כך למדוד באופן יחסי באופן הדף שבו אדם אחד אדם אחד
 אדם אחד.

2) דיווחים (אדם אחד שבו אדם אחד)

1, 3, 4, 8, 9, 10, 13, 14, 18

Recall נקודה	precision	הנקודה	Recall	נקודה
0	1/1 = 1			
0.1	3/4 = 0.75			
0.2	4/8 = 0.5			
0.3	6/10 = 0.6			
0.4	7/13 = 0.538			
0.5	8/14 = 0.57			
0.6	0			
...	...			
1	0			

$$\frac{1 + 0.75 + 0.5 + 0.3 + 0.538 + 0.57}{11} = 0.3599$$

$$P_1 - \text{precision} = \frac{3}{6} = 0.5$$

3 (2)

$$\text{precision}_{e10} = \frac{6}{10}$$

7

NPNG ES number ~~number~~ number 1
 Recall " number 2

3

2

72

$$W_{\text{no Rev, 3}} = \frac{2}{3} \cdot \frac{5}{4} = \frac{10}{12}$$

$$W_{\text{mother}, 4} = \frac{1}{3} \cdot \frac{5}{4} = \frac{5}{12}$$

$$W_m \quad 15 = \frac{3}{3} \cdot \frac{5}{4} = \frac{5}{4}$$

$$WF_{\text{after}, 2} = \frac{4}{4} \cdot \frac{5}{3} = \frac{5}{3}$$

$$WF_{13} = \frac{1}{3} \cdot \frac{5}{3} = \frac{5}{9}$$

$$W_{F,5} = \frac{2}{3} \cdot \frac{5}{3} = \frac{10}{9}$$

$$\text{Sim}(\text{Dir } q) = \text{Inner product} = \sum w_{ij} \cdot q_{ij}$$

$$\sin(\theta_{11q}) = \frac{15}{16} = 0.9375$$

$$\text{Sim}(D_2, q) = \frac{5}{3} = 1.66$$

$$S.m(D3, q) = \frac{10}{12} + \frac{5}{q} = 1.38$$

$$\text{Sim}(D_4, q) = \frac{r}{12} = 0.416$$

$$\sin(15,9) = \frac{r}{9} = \frac{10}{9} = 2.36$$

ה'תשס"ח י"ח שבט ה'תשס"ח

$$D_5 \rightarrow D_2 \rightarrow D_3 \rightarrow D_1 \rightarrow D_4$$

(2) היתכנות לא ישתנה, כיוון שכך הכנסה בסף, יין שלוש שנים, צבר שיעור לא מספיק הצירוף, יין לא ילפף אלא תחילה.

(3) מלפני 2 צברים יקרא: (1) יענה ה f זה u $(\frac{5}{3})$

(2) ה f זה u 5 יענה.

סך נילאקאלה אלא מחשבים עם 2, 3, 5 יענה. האשנים
אם כן לא, קולם, אכן השפלה. (כאן 5 יענה עם 2, 3 יענה השפלה)

(4) סך $k=2$ $f(\vec{x}, \vec{y}) = |x_1 - y_1| + |x_1 - x_2|$

cluster 1
 d_1

cluster 2
 d_2, d_3

$$f(d_1, d_3) = |0.8 - 0.4| + |0.4 - 0.2| = 1.1$$

$$f(d_2, d_3) = 0.9$$

$$f(d_1, d_4) = 1$$

$$f(d_1, d_5) = 1$$

$$f(d_2, d_5) = 0.8$$

$$f(d_1, d_6) = 1$$

$$f(d_2, d_6) = 0.7$$

$$f(d_1, d_7) = 0.8$$

$$f(d_2, d_7) = 0.6$$

$$c_1 = (0.9, 0.8)$$

$$c_2 = \frac{1}{6} (d_2 + d_3 + d_4 + d_5 + d_6 + d_7) = \frac{1}{6} (0.5 \quad 0.3)$$

cluster 1
d₁
d₂

cluster 2
d₃
d₄
d₅
d₆
d₇

$$f(d_1, c_2) = 0.88$$

$$f(d_1, c_1) = 0$$

$$f(d_2, c_1) = 0.2$$

$$f(d_2, c_2) = 0.68$$

$$f(d_3, c_1) = 1.1$$

$$f(d_3, c_2) = 0.9$$

$$f(d_4, c_1) = 1.1$$

$$f(d_4, c_2) = 0.4$$

$$f(d_5, c_1) = 1$$

$$f(d_5, c_2) = 0.12$$

$$f(d_6, c_1) = 1$$

$$f(d_6, c_2) = 0.12$$

$$f(d_7, c_1) = 0.8$$

$$f(d_7, c_2) = 0.1$$

(4)

70%

$$c_1 = \frac{1}{2} [d_1 + d_2] = (0.75 \ 0.75)$$

$$c_2 = \frac{1}{5} [d_3 + d_4 + d_5 + d_6 + d_7] = (0.46 \ 0.24)$$

cluster 1 cluster 2

d₁
d₂

d₃
d₄
d₅
d₆
d₇

$$f(d_1, c_1) = 1$$

$$f(d_1, c_2) = 1$$

$$f(d_2, c_1) = 0.1$$

$$f(d_2, c_2) = 0.8$$

$$f(d_3, c_2) = 0.42$$

$$f(d_4, c_1) = 1$$

$$f(d_4, c_2) = 0.38$$

$$f(d_5, c_1) = 0.9$$

$$f(d_5, c_2) = 0.08$$

$$f(d_6, c_1) = 0.9$$

$$f(d_6, c_2) = 0.2$$

$$f(d_7, c_1) = 0.7$$

$$f(d_7, c_2) = 0.2$$

55% 70% 80%

11

70%

==

4) משימה זו היא - באיזה מיון אנו יורד מהר (אנחנו חוקרים משימה)

1) (b) clustering - זה יוצר מרכזים אלו הקטגוריות והם קטגוריות
זה אקראיות אלו המשימות.

2) (a) categorizing - יוצר מרכזים אלו הקטגוריות אין לסווגם
אלו המשימות.

3) (c) feature - זה יוצר משימה עם clustering על משימות
המשימה היא "למצוא" את המשימה וזהו המשימה.

feature selection - זה יוצר משימה, ואנחנו חוקרים משימה
היא אקראית - אין אנו חוקרים משימה.

