1至Q&A 😌

이정인, 신민경, 전성원, 김덕영

Q. 수업시간에 다뤘던 알고리즘은 오버랩 카운트를 재사용하지 않았는데, PPJoin이 맞는지 궁금합니다.

```
for i in range(len(R)): # 후보 검증 후 S에 추가 (y_id, y) = R[i]

if A[i] == 0: continue # prefix 안에 공통 토큰이 하나도 없다면 제외
a = math.ceil((len(x)+len(y))*r/(r+1))
```

검증 코드 부분에서 count = 0 이면, 공통 토큰을 더이상 count하지 않는 방법으로 재사용합니다.

저희는 최종적으로 가장 효율적인 알고리즘을 APJoin로 선정해서 구현한 것이므로 PPJoin 방식과는 다릅니다.

Q. Candidate Pair을 줄이는 과정에서, 실제 threshold를 넘을 후보 pair의 손실은 없었나요?

[수도 코드]

```
11 for i = 1 to max_probe_prefix do
12 w ← x[i]; // 레코드 x의 토큰 (x, i)
13 for each (y, j) ∈ I<sub>w</sub> do // 레코드 y의 토큰 (y, j)
14 if |y| < t ⋅ |x| // 현재 노드 이후 삭제
15 (y, j)와 연결된 노드들 삭제;
16 else if j > prefix_y[|x|-|y|+1]
17 (y, j) 삭제; // 현재 노드 삭제
18 else if i > prefix_x[|x|-|y|+1]
19 continue; // 다음 노드로 이동
20 else
21 A[y] = A[y] + 1; // candidate generation
```

$$O(x,y) \ge \alpha = \frac{t}{1+t}(|x|+|y|)$$

$$\alpha \le |y|, \quad \frac{t}{1+t}(|x|+|y|) \le |y|$$

$$t|x| \le |y|$$

Length filtering과 prefix filtering 에서는

후보가 될 가능성이 전혀 없는 경우만 제외하기 때문에 후보 Pair에 대한 손실이 없습니다.

Q. 혹시 직접 생각한 알고리즘이 있나요?

```
for i=1 to max_index_prefix do w \leftarrow x[i]; I_w = I_w \cup \{(x,\ i)\}; // 토큰 (x,\ i)를 inverted index에 VerifyZip(x,\ A,\ \alpha); // verification phase
```

```
for i in range(len(R)): # 후보 검증 후 S에 추가
   (y_id, y) = R[i]
   if A[i] = 0: continue # prefix 안에 공통 토큰이 하나도 없다면 제외
   a = math.ceil((len(x)+len(y))*r/(r+1))
   (p,q,cnt) = (0,0,0)
   while(p<len(x) and q<len(y)):# 공통 토큰 개수 카운트
       if x[p] = v[q]:
          cnt = cnt+1
          p = p+1
          a = a+1
       elif cnt+min(len(x)-p-1,len(y)-q-1) < a: break # positional filtering
       elif x[p] < y[q]: p = p+1
      else: q = q+1
   if cnt >= a and y_id not in x: # x와 y가 이미 친구인 경우는 제외
       S.append((x id. v id))
```

논문에서 제시된 APJoin 알고리즘은 검증에 관한 VerifyZip 코드에 대한 설명과 구현이 안되어 있었는데,

- 1) 이에 대해 저희는 VeritifyZip 코드를 직접 구현하였습니다.
- 2) 공통토큰을 세기 위해 O(n^2) 방법 대신 O(n) 방법으로 개선하는 방향으로 코드를 생각하여 구현했습니다.
- 3) 이미 친구인 경우에 어떻게 후보군을 제외시킬 지에 대해 직접 고민해봤습니다.

Q. 최종 결과가 어떻게 나왔는지 확인 가능할까요?

```
1 1 2
2 3
3 2 5
4 4 3
5 6 3
6 4 6
```

```
# when r = 0.1
[(5, 1), (2, 6), (2, 4), (3, 1), (3, 5)]
duration: 0.0010955333709716797
# when r = 0.6
[(5, 1)]
duration: 0.0004329681396484375
# when r = 0.7
[(5, 1)]
duration: 0.0004551410675048828
# when r = 0.8
[(5, 1)]
duration: 0.0003483295440673828
# when r = 0.9
[(5, 1)]
duration: 0.0003559589385986328
```

Simple Test Case 를 생성하여 결과를 돌려봤습니다.

Q. 자료조사 열심히 하신 것 잘 봤습니다. 그 많은 알고리즘들을 다 적용해 보신건가요?

논문 내 4가지 Join들의 성능비교

	Comp	1) Cand	Join	Time(sec)
PPJoin 2)	4,132,651	1,356,337	3,683	2.106
PPJoin+	4,132,651	11,905	3,683	1.810
MPJoin	2,200,083	1,356,628	3,683	0.978
APJoin	1,378,845	1,357,015	3,683	0.668

- 1) PPJoin vs. PPJoin+: 비교연산 수는 동일, 추출된 조인후보 쌍은 큰 차이
- 2) PPJoin vs. (MPJoin & APJoin): 비교연산 수는 큰 차이, 추출된 조인후보 쌍은 차이 없음

BruteForce 와 수업시간에 간단한 Similarity Join 관련 알고리즘에 대해 구현을 해보며 이해를 했었고, PPJoin과 PPJoin+, MPJoin 알고리즘에 대한 수도코드 분석을 통해 논문을 이해하고 정리했습니다.

논문에서도 성능검증이 완료된 표를 통해 알 수 있었듯이 최종적으로 가장 효율적인 알고리즘이 'APJoin' 을 최종 구현했습니다.

Q. 문제에서 해결해야 하는 친구 추천 코드 구현 시, pyspark는 어떤 식으로 이용했나요?

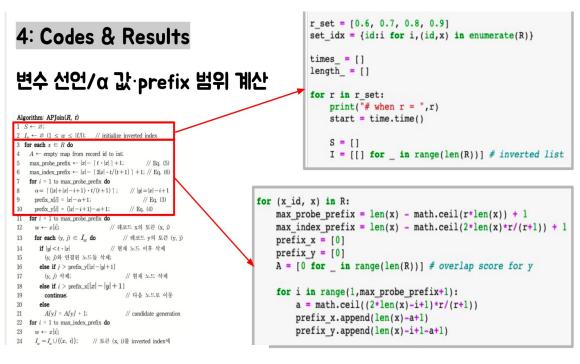
```
In [1]: import findspark
        findspark.init()
In [2]: from pyspark import SparkContext
        from pyspark.sql import SQLContext, SparkSession
        sc = SparkContext("local", "pbl-1")
        spark = SparkSession(sc)
In [3]: data = sc.textFile("facebook combined.txt")
        parse data = data.flatMap(lambda line: line.split('\n')) \
                    .map(lambda l:l.split(' ')).map(lambda l: [(int(l[0]),int(l[1])),(int(l[1]),int(l[0]))]).flatMap(lambda t:
In [4]: group = parse data.groupByKey().mapValues(list)
        new group = group.flatMap(lambda 1: [(l[0],(len(l[1]), sorted(l[1])))])
        #print(new group.collect())
In [5]: order = new group.sortBy(lambda v: v[1])
        R = order.mapValues(lambda v: v[1]).collect()
        #print(R)
```

Q. 문제에서 해결해야 하는 친구 추천 코드 구현 시, pyspark는 어떤 식으로 이용했나요?

4: Codes & Results

동적 Prefix Filtering/Inverted list 작성

```
for i in range(1, max probe prefix+1):
Algorithm: APJoin(R, t)
1 S ← Ø;
                                                                                w = x[i-1]
2 I_w \leftarrow \emptyset (1 \leq w \leq |U|); // initialize inverted index
3 for each x \in R do
                                                                                w idx = set idx[w]
4 A ← empty map from record id to int;
   max probe prefix \leftarrow |x| - \lceil t \cdot |x| \rceil + 1:
   max_index_prefix \leftarrow |x| - \lceil 2|x| \cdot t/(t+1) \rceil + 1; // Eq. (6)
                                                                                for ((y id, y), j) in I[w idx]:
   for i = 1 to max_probe_prefix do
                                                                                      if len(y) < r*len(x):</pre>
    \alpha = \lceil (|x|+|x|-i+1) \cdot t/(t+1) \rceil; // |y| = |x|-i+1
                                                                                            I[w idx].remove(((y id,y),j))
     prefix_x[i] = |x| - \alpha + 1;
                                     // Eq. (3)
    prefix v[i] = (|x|-i+1)-\alpha+1;
                                 // Eq. (4)
                                                                                      elif j > prefix y[len(x)-len(y)+1]:
11 for i = 1 to max probe prefix do
                                                                                            I[w_idx].remove(((y_id,y),j))
                           // 레코드 x의 토큰 (x, i)
                            // 레코드 y의 토큰 (y, j)
     for each (y, j) \in I_{m} do
                                                                                      elif i > prefix x[len(x)-len(y)+1]: continue
                           // 현재 노드 이후 삭제
       (y, j)와 연결된 노드들 삭제;
       else if j > \text{prefix}_y[|x| - |y| + 1]
                                                                                            A[set idx[y id]] = A[set idx[y id]] + 1 # 후보 확정
       (y, j) 삭제;
       else if i > \text{prefix}_x[|x| - |y| + 1]
                                // 다음 노드로 이동
                                                            for i in range(1,max index prefix+1): # x의 원소에 대한 inverted list 추가
        A[y] = A[y] + 1;
                                // candidate generation
22 for i = 1 to max_index_prefix do
23 w ← x[i];
                                                                  w idx = set idx[w]
    I_{w} = I_{w} \cup \{(x, i)\};
                     // 토큰 (x, i)를 inverted index에
                                                                  I[w idx].append(((x id,x),i))
    VerifyZip(x, A, \alpha);
26 return S
                                                                                   조인 후보 쌍 검증
```



```
Algorithm: APIoin(R. t)
                                                                               for i in range(len(R)): # 후보 검증 후 S에 추가
1 S ← Ø;
                                                                                     (y id, y) = R[i]

 I. ← Ø (1 < w < |II|); // initialize inverted index</li>

3 for each x \in R do
4 A ← empty map from record id to int;
                                                                                     if A[i] == 0: continue # prefix 안에 공통 토큰이 하나도 없다면 제외
5 max_probe_prefix \leftarrow |x| - \lceil t \cdot |x| \rceil + 1;
   max_index_prefix \leftarrow |x| - \lceil 2|x| \cdot t/(t+1) \rceil + 1; // Eq. (6)
                                                                                     a = math.ceil((len(x)+len(y))*r/(r+1))
   for i = 1 to max probe prefix do
     prefix x[i] = |x| - \alpha + 1;
                                     // Eq. (3)
                                                                                     (p,q,cnt) = (0,0,0)
    prefix_v[i] = (|x|-i+1)-\alpha+1;
11 for i = 1 to max_probe_prefix do
                            // 레코드 x의 토큰 (x, i)
    for each (y, j) \in I_w do
                              // 레코드 v의 토큰 (v. j)
                                                                  while(p<len(x) and q<len(y)):# 공통 토큰 개수 카운트
                            // 현재 노드 이후 삭제
      if |y| < t \cdot |x|
                                                                       if x[p] == y[q]:
        (y, j)와 연결된 노드들 삭제;
                                                                             cnt = cnt+1
       else if j > \text{prefix}_y[|x| - |y| + 1]
       (y, j) 삭제;
                                // 현재 노드 삭제
                                                                             p = p+1
       else if i > \operatorname{prefix}_{\mathbf{x}}[|x| - |y| + 1]
                                                                             q = q+1
                                 // 다음 노드로 이동
                                                                        elif cnt+min(len(x)-p-1,len(y)-q-1) < a: break # positional filtering
                                                                       elif x[p] < y[q]: p = p+1
        A[y] = A[y] + 1;
                                 // candidate generation
22 for i = 1 to max index prefix do
                                                                        else: q = q+1
                                                                  if cnt >= a and y id not in x: # x와 y가 이미 친구인 경우는 제외
                                                                       S.append((x_id, y_id))
25 VerifyZip(x, A, α);
```

Q. 문제에서 해결해야 하는 친구 추천 코드 구현 시, pyspark는 어떤 식으로 이용했나요?

RDD로 구현한 코드

```
# when r = 0.7
for (x_id, x) in R: # 순서대로 inverted list를 작성해야 하므로 동기
   A = [0 \text{ for } \_ \text{ in } range(len(R))]
   max\_probe\_prefix = len(x) - math.ceil(r*len(x)) + 1
   max_index_prefix = len(x) - math.ceil(2*len(x)*r/(r+1)) + 1
   prefix_x = sc.parallelize([i for i in range(max_probe_prefix+1)]) \( \psi
  \)
                .map(lambda t: len(x)-(math.ceil((2*len(x)-t+1)*r/(r+1)))+1).collect()
   prefix_y = sc.parallelize([i for i in range(max_probe_prefix+1)]) \( \psi
  \)
                .map(lambda t: len(x)-t+1-(math.ceil((2*len(x)-t+1)*r/(r+1)))+1).collect()
    for i in range(1,max_probe_prefix+1):
       w = x[i-1]
       w_{idx} = set_{idx}[w]
        inverted_list = sc.parallelize([w_idx]).filter(lambda t: len(t[0][1]) >= r*len(x) and t[1] <= prefix_y[len(x)-len(t[0][1])+1])
        [[w idx] = inverted list.collect()
       candidate = inverted_list.filter(lambda t: i \leftarrow prefix_x[len(x)-len(t[0][1])+1]).collect()
        for item in candidate:
            item_id = item[0][0]
            A[set idx[item id]] = A[set idx[item id]]+1
```

duration: 699.9134256839752

Q. Positional filtering에서 후보군 제거 과정에서, 어떻게 제거한 건지 궁금합니다.

3: Algorithm Research - Positional Filtering

: Overlap의 최대크기 ≥ α 를 만족하지 못하는 레코드를 조인후보 쌍 생성 전에 제외한다.

레코드 x, y의 공통토큰이 각각 i번째, j번째 있다고 하자.

Overlap의 최대크기 = 현재 Overlap 크기 + min(len(x)-i, len(y)-j)

prefix filtering은 prefix 범위 내에서는

O(x,y) ≥ α를 만족하지 못하더라도 조인후보 쌍으로 선택한다.

-> 불필요한 조인후보 쌍이 많이 생길 수 있다.

∴ positional filtering 이용하면, 생성되는 조인후보의 쌍을 많이 줄일 수 있다.