

Computer Convergence Application

Homework 3: Frequent pattern mining

컴퓨터공학부

2017-18538

황선영

1. finding candidate motif sequence

Drug-induced liver injury(DILI)는 가장 예측하기 어려운 항생제에 대한 부작용 중 하나이며 승인된 의약품이 시판되고도 철회하게 하는 주요 원인이다. FDA는 약 1000여개에 대한 약물에 대해 DILI severity와 potential 에 따라 분류하여 dataset을 구축하였다. Dataset에 있는 약물에 대해 높은 확률로 DILI를 유발하는 약물과 DILI를 유발하지 않는 약물 각각에 대해 pattern mining algorithm을 이용하여 candidate motif sequence를 찾음으로써 어떤 motif sequence가 DILI 유발 여부에 관여하는지 분석한다.

vDILICConcern인 vMost_DILI-Concern인 compound, 즉 거의 drug-induced liver injury를 causing하는 약물을 positive로 정한 다음 frequent하게 나타나는 motif sequence를 찾는다. 이 때 negative에 있는 motif sequence는 제외하기 때문에 positive에만 있는 motif sequence만 pure_positive_result로 나타난다. vNo-DILI-Concern인 compound, 즉 drug-induced liver injury를 causing하지 않는 약물을 negative로 정한 다음 frequent하게 나타나는 motif sequence를 찾는다. 이 때 positive에 있는 motif sequence는 제외하기 때문에 negative에만 있는 motif sequence만 pure_negative_result로 나타난다. share_positive_result와 share_negative_result는 positive와 negative에 동시에 나타나는 candidate motif sequence를 의미한다. 두 결과가 같으므로 share_positive_result만 첨부하였다. 그리고 data의 양이 많아 index 0-4에 해당하는 요소 5개만 나타냈다.

2. pattern mining algorithm: Apriori vs FP growth

Apriori algorithm은 apriori property와 join과 prune property를 이용하는 pattern mining algorithm이다.

```
126
[(['2cc'], 0.25955414012738853),
 (['O)c'], 0.29904458598726114),
 (['cc2'], 0.32961783439490444),
 (['c(c', '(=O'], 0.2764331210191083),
 (['c1c', '(=O'], 0.3200636942675159),
```

[Figure1-1] Apriori algorithm: pure_positive_result

```
29
[(['(O)'], 0.35940860215053766),
 ([')O'], 0.2615591397849462),
 (['1C'], 0.25887096774193546),
 (['CC('], 0.3293010752688172),
 (['CCC'], 0.3682795698924731),
```

[Figure1-2] Apriori algorithm: pure_negative_result

```

57
[(['(=O'], 0.5022292993630574),
(['(C('], 0.26528662420382165),
(['(C)'], 0.40732484076433123),
(['(CC'], 0.27038216560509554),
(['(cc'], 0.38407643312101913),

```

[Figure1-3] Apriori algorithm: share_positive_result

FP-growth algorithm은 minimum support를 만족하는 database로부터 conditional frequent pattern tree와 conditional pattern을 construct하는 pattern mining algorithm이다. Scan 횟수가 apriori algorithm에 비해 적어 execution time이 짧다. Execution time이 중요하다면 FP-growth algorithm을 선택하면 좋을 것이다.

```

144
[(['O)c'], 0.2984076433121019),
(['cc2'], 0.3305732484076433),
(['2cc'], 0.25127388535031847),
(['c1c', 'c2c'], 0.2843949044585987),
(['c2c', '=O)'], 0.30668789808917196),

```

[Figure2-1] FP-growth algorithm: pure_positive_result

```

27
[(['CCN'], 0.29811827956989245),
(['1)C'], 0.27338709677419354),
(['CC('], 0.32016129032258067),
(['CCC'], 0.35940860215053766),
(['ccc', '(CC'], 0.29596774193548386),

```

[Figure2-2] FP-growth algorithm: pure_negative_result

```

54
[(['c2c'], 0.439171974522293),
(['c(c'], 0.5226114649681529),
(['c1c'], 0.6095541401273885),
(['c1)'], 0.4194267515923567),
(['=O)'], 0.6917197452229299),

```

[Figure2-3] FP-growth algorithm: share_positive_result

3. length of patterns(n-gram)

n-gram은 연속적인 n개의 token(word, character)로 구성된 것이다. 문장 분석의 경우 n-gram이 커질수록 맥락을 잘 내포하는 경향이 있다. N-gram이 3일 때의 결과는 Figure1과 같다. N-gram이 2인 경우 pure_positive의 길이가 5000이 넘을 만큼 많아서 결과가 유효하지 않다고 판단해 첨부하지 않았다. N-gram이 4일 때와 5일 때를 확인해보니 candidate motif sequence의 구조가 더 명확해질 수 있겠으나 그 수가 굉장히 적은 것을 알 수 있었다. N-gram이 3일 때 적절한 수의

candidate motif sequence가 도출되므로 n-gram=3이 가장 적절하다고 판단된다.

```
11
[(['(=O)'], 0.5028662420382166),
([')=O)'], 0.3375796178343949),
(['1ccc'], 0.34394904458598724),
(['C(=O)'], 0.36528662420382163),
(['c(cc)'], 0.36528662420382163),
```

[Figure3-1] n-gram=4: pure_positive_result

```
4
[(['(C)C'], 0.28440860215053765),
(['=O)C'], 0.2801075268817204),
(['ccc1'], 0.29731182795698924),
(['cccc', 'c1cc'], 0.2629032258064516)]
```

[Figure3-2] n-gram=4: pure_negative_result

```
2
[(['cc(c)'], 0.3213375796178344), ([ 'c1cc', 'ccc('],
0.2535031847133758)]
```

[Figure3-3] n-gram=4: share_positive_result

```
0
[ ]
```

[Figure4-1] n-gram=5: pure_positive_result

```
1
[(['cccc'], 0.2577956989247312)]
```

[Figure4-2] n-gram=5: pure_negative_result

```
2
[(['C(=O)'], 0.36942675159235666), ([ 'c1ccc'], 0.31751592356687897)]
```

[Figure4-3] n-gram=5: share_positive_result

4. minimum support

Minimum support가 커질 수록 support가 높아지는 대신 candidate motif sequence의 수가 급격하게 적어지는 것을 알 수 있다. 하지만 minimum support가 너무 낮다면 support가 낮은 candidate motif sequence가 pure_positive_result, pure_negative_result에 나타나게 되는데 support가 낮으므로 그 집단의 대표성을 띄기 힘들다. Minimum support가 0.25일 때의 결과는 Figure1과 같다. Minimum support가 0.25-0.375사이의 범위에 있어야 support도 적당히 높고 대신 candidate motif sequence의 수도 적당한 결과가 도출될 것이다.

```

5
[(['(CC'], 0.43172043010752686),
 (['])CC'], 0.3879032258064516),
 ([ 'C(C'], 0.435752688172043),
 ([ 'C)C'], 0.46344086021505376),
 ([ 'O)C'], 0.4344086021505376)]

```

[Figure5-1] minimum_sup=0.375: pure_positive_result

```

14
[(['(=O'], 0.5),
 ([ ' (C)'], 0.4105095541401274),
 ([ ')=O'], 0.478343949044586),
 ([ '1cc'], 0.4621019108280255),
 ([ '=O)'], 0.6917197452229299),

```

[Figure5-2] minimum_sup=0.375: pure_negative_result

```

22
[(['(cc'], 0.3945859872611465),
 ([ ' )cc'], 0.38980891719745225),
 ([ 'C(='], 0.39363057324840767),
 ([ 'c(c'], 0.5226114649681529),
 ([ 'c1)'], 0.4194267515923567),

```

[Figure5-3] minimum_sup=0.375: share_positive_result

```

7
[(['(=O'], 0.5),
 ([ 'c(c'], 0.5226114649681529),
 ([ 'cc('], 0.5837579617834395),
 ([ 'cc1'], 0.5089171974522293),
 ([ '(=O', '=O)'], 0.5),

```

[Figure5-1] minimum_sup=0.5: pure_positive_result

```

0
[ ]

```

[Figure5-2] minimum_sup=0.5: pure_negative_result

```

3
[(['=O)'], 0.6917197452229299),
 ([ 'c1c'], 0.6095541401273885),
 ([ 'ccc'], 0.6910828025477707)]

```

[Figure5-3] minimum_sup=0.5: share_positive_result

5. augmentation & tokenization

Augmentation은 data에 인위적인 변화를 가해 새로운 training data를 대량 확보하는 기법이다. 현실

세계에서도 실제로 존재할 법한 데이터를 생성함으로써 좀 더 일반화된 model을 얻는 것을 목표로 한다. Augmentation이 20일 때의 결과는 Figure1과 같다. Augmentation은 기존의 data의 정보량은 보존한 상태로 noise를 주는 방식이기 때문에 정보량 자체는 변하지 않는다. 정보량에 약간의 변화를 주는 것이기 때문에 deep learning으로 분석된 data의 표현되는 강력한 특징을 느슨하게 만드는 것이다. 이는 overfitting을 막아주고 예측 범위를 약간 넓혀줄 수는 있으나 매우 큰 효과를 가져오진 못한다.

```
65
[(['(=O'], 0.40483870967741936),
(['(C('], 0.3317204301075269),
(['(C)'], 0.47096774193548385),
(['(CC'], 0.42258064516129035),
(['(cc'], 0.28440860215053765),
```

[Figure6-1] Augmentation=10: pure_positive_result

```
21
[(['(O)'], 0.3639784946236559),
([')O)'], 0.2564516129032258),
(['CC('], 0.32150537634408605),
(['CCC'], 0.3548387096774194),
(['CCN'], 0.30053763440860215),
```

[Figure6-2] Augmentation=10: pure_negative_result

```
157
[(['O)c'], 0.29044585987261146),
(['cc2'], 0.321656050955414),
(['(=O', '1cc'], 0.2681528662420382),
(['c(c', '(=O'], 0.29745222929936305),
(['c1c', '(=O'], 0.3452229299363057),
```

[Figure6-3] Augmentation=10: share_positive_result