

# Assignment 2: ASQL, ICs, JDI

Student: Koray Poyraz

Studentnr: 5367840

## Question 1

see file db.sql

In the creation of the relations I used primary keys, foreign keys, cascading, set null, domain types, not nulls, check ...

### Department

Buildings is not unique because it might be that a building might have multiple departments, so building could have his own relation but due to not having a specific attributes related to a building its not needed.

Student and instructor

These relations have a null on either side of the relation because we don't want to delete tuples from student and instructor when a department gets deleted, it might be that a student can be an instructor and join another department.

Advisor

If a student is deleted or updated then those tuples in the advisor are deleted and updated as well. If an instructor is deleted then we set the value to null so you can have in a glance all students who don't have a advisor. Also, it has cascade on update.

Course

I did not make title a unique attribute because it might be that a department wants to give the same course with same title but a different course id. We could make title and dept name unique but that would not give the possibility to have the same course with different id, same title within the same department but different credits.

When a course is deleted or updated then the effect will be cascaded to the tuples in prereq relation.

Teaches

When we remove a section then there is no need to have the related teaches tuples.

Section

If there is a new classroom for a section then the update will reflect the change or delete the related tuples on delete. If a course is been deleted then there is no need to have the related tuples in section.

Takes

When a student is deleted or updated the change will be reflected on the related tuples in takes. If a section is deleted then the related tuples in takes will be deleted as well and on update the change will be cascaded.

## Question 2

see file inserts.sql

## Question 3

R(A,B,C,D,E,F,G)

Functional dependencies:

- ABB -> EG
  - C -> DG
  - E -> FG
  - AB -> C
  - G -> F
- Closure
- ABCDEFG+ -> {ABCDEFG}
  - ABD can determine EG so we can discard EG
  - and because E and G determines F we can also discard F (Transitivity)
  - ABCD+ -> {ABCDEFG}
  - C can determine DG so we can discard DG
  - ABC+ -> {ABCDEG}
  - AB can determine C so we can discard C
  - AB+ -> {ABCDEG}

The proper subset of AB+ is (A),(B), A and B are not super key

So below I have a code which shows the first name -> genders values which are violated, total is 52 first names

# the code could be more simplified by just using the unique() and intersection() methods

# but I found it important to show both column values

violated values = df.name.gender\

.groupby('first\_name')\

.apply(lambda x: ~ ' '.join(x['gender']))\

.reset\_index()\

.rename(columns={'count\_diff\_gender':''})

count\_names = 0

for row in violated\_values.values:

set\_val = set(row[1].split(','))

if len(set\_val) > 1:

count\_names += 1

# print('First names: [row[0]] \t Gender: {set\_val}')\

# print('Violated: (count\_names)')

# Version 2: shows all tuples including the row number, the name and gender which violate

key\_map = {}

for index, x in enumerate(df.name\_gender.values):

if x[0] not in key\_map:

key\_map[x[0]] = [index, x[1]]

else:

idx\_of\_name = key\_map[x[0]]

is\_violated = len(set(idx\_of\_name[1])) > 1

count = 0

for name in key\_map:

genders = map(lambda x: x[1], key\_map[name])

is\_violated = len(set(genders)) > 1

if is\_violated:

for gender in key\_map[name]:

gender\_idx = gender.split(',')

count += 1

print(f'Row number: {row[0]} \t name: {name} \t gender: {gender\_idx[1]}')

print(f'Violated: {count\_names}')

print(f'Number of (unique) first names which violate: {count\_names}')

print(f'Number of tuples which violate: {count}')

Violations:

Row number: [1] name: David gender: M

Row number: [144] name: David gender: M

Row number: [187] name: David gender: M

Row number: [201] name: David gender: M

Row number: [263] name: David gender: F

Row number: [203] name: David gender: M

Row number: [300] name: David gender: M

Row number: [284] name: David gender: M

Row number: [346] name: David gender: M

Row number: [377] name: David gender: M

Row number: [545] name: David gender: M

Row number: [460] name: David gender: M

Row number: [546] name: David gender: M

Row number: [678] name: David gender: M

Row number: [652] name: David gender: M

Row number: [633] name: David gender: M

Row number: [688] name: David gender: M

Row number: [799] name: David gender: M

Row number: [845] name: David gender: M

Row number: [859] name: David gender: M

Row number: [881] name: David gender: M

Row number: [917] name: David gender: M

Row number: [968] name: David gender: M

Row number: [969] name: David gender: M

Row number: [1086] name: David gender: M

Row number: [1072] name: David gender: M

Row number: [1289] name: David gender: M

Row number: [1290] name: David gender: M

Row number: [1318] name: David gender: M

Row number: [1322] name: David gender: M

Row number: [1347] name: David gender: M

Row number: [1349] name: David gender: M

Row number: [1626] name: David gender: M

Row number: [1673] name: David gender: M

Row number: [1693] name: David gender: M

Row number: [1702] name: David gender: M

Row number: [1761] name: David gender: M

Row number: [1762] name: David gender: M

Row number: [1896] name: David gender: M

Row number: [1958] name: David gender: M

Row number: [1991] name: David gender: M

Row number: [2007] name: David gender: M

Row number: [2021] name: David gender: M

Row number: [2187] name: David gender: M

Row number: [2192] name: David gender: M

Row number: [2130] name: David gender: M

Row number: [2173] name: David gender: M

Row number: [2232] name: David gender: M

Row number: [2253] name: David gender: M

Row number: [2317] name: David gender: M

Row number: [2392] name: David gender: M

Row number: [2603] name: David gender: M

Row number: [2649] name: David gender: M

Row number: [2755] name: David gender: M

Row number: [2940] name: David gender: M

Row number: [3013] name: David gender: M

Row number: [3396] name: David gender: M

Row number: [3103] name: David gender: M

Row number: [3119] name: David gender: M

Row number: [3235] name: David gender: M

Row number: [3240] name: David gender: M

Row number: [3355] name: David gender: M

Row number: [3597] name: David gender: M

Row number: [3584] name: David gender: M

Row number: [3627] name: David gender: M

Row number: [3640] name: David gender: M

Row number: [3681] name: David gender: M

Row number: [3688] name: David gender: M

Row number: [3805] name: David gender: M

Row number: [3967] name: David gender: M

Row number: [4077] name: David gender: M

Row number: [4185] name: David gender: M

Row number: [4235] name: David gender: M

Row number: [4238] name: David gender: M

Row number: [4293] name: David gender: M

Row number: [4356] name: David gender: M

Row number: [4453] name: David gender: F

Row number: [4517] name: David gender: M

Row number: [4592] name: David gender: M

Row number: [4606] name: David gender: M

Row number: [4618] name: David gender: M

Row number: [4883] name: David gender: M

Row number: [4734] name: David gender: M

Row number: [4769] name: David gender: M

Row number: [4795] name: David gender: M

Row number: [4962] name: David gender: M

Row number: [5033] name: David gender: M

Row number: [5094] name: David gender: M

Row number: [5141] name: David gender: M

Row number: [5221] name: David gender: M

Row number: [5361] name: David gender: M

Row number: [5476] name: David gender: M

Row number: [5526] name: David gender: M

Row number: [5691] name: David gender: M

Row number: [5792] name: David gender: M

Row number: [5893] name: David gender: M

Row number: [6019] name: David gender: M

Row number: [6042] name: David gender: M

Row number: [6193] name: David gender: M

Row number: [6207] name: David gender: M

Row number: [6269] name: David gender: M

Row number: [6276] name: David gender: M

Row number: [6355] name: David gender: M

Row number: [6387] name: David gender: M

Row number: [6098] name: David gender: M

Row number: [6408] name: David gender: M

Row number: [6480] name: David gender: M

Row number: [6535] name: David gender: M

Row number: [6541] name: David gender: M

Row number: [6551] name: David gender: M

Row number: [7847] name: David gender: M

Row number: [6560] name: David gender: M

Row number: [6574] name: David gender: M

Row number: [6604] name: David gender: M

Row number: [6619] name: David gender: M

Row number: [6756] name: David gender: M

Row number: [6821] name: David gender: M

Row number: [6862] name: David gender: M

Row number: [6897] name: David gender: M

Row number: [6965] name: David gender: M

Row number: [7073] name: David gender: M

Row number: [7087] name: David gender: M

Row number: [7104] name: David gender: M

Row number: [7111] name: David gender: M

Row number: [7231] name: David gender: M

Row number: [7260] name: David gender: M

Row number: [7302] name: David gender: M

Row number: [7350] name: David gender: M

Row number: [7380] name: David gender: M

Row number: [7394] name: David gender: M

Row number: [7456] name: David gender: M

Row number: [7558] name: David gender: M

Row number: [7567] name: David gender: M

Row number: [7644] name: David gender: M

Row number: [7622] name: David gender: M

Row number: [7840] name: David gender: M

Row number: [7862] name: David gender: M

Row number: [7909] name: David gender: M

Row number: [7997] name: David gender: M

Row number: [8006] name: David gender: M

Row number: [8039] name: David gender: M

Row number: [8409] name: David gender: M

Row number: [8512] name: David gender: M

Row number: [8541] name: David gender: M

Row number: [8548] name: David gender: M

Row number: [8634] name: David gender: M

Row number: [8683] name: David gender: M

Row number: [8779] name: David gender: M

Row number: [8798] name: David gender: M

Row number: [8943] name: David gender: M

Row number: [8987] name: David gender: M

Row number: [8998] name: David gender: M

Row number: [9010] name: David gender: M

Row number: [9098] name: David gender: M

-----

Row number: [86] name: Jean gender: M

Row number: [291] name: Jean gender: F

Row number: [1169] name: Jean gender: F

Row number: [1170] name: Jean gender: F

Row number: [317] name: Jean gender: F

Row number: [432] name: Jean gender: M

Row number: [4624] name: Jean gender: M

Row number: [5359] name: Jean gender: F

Row number: [5913] name: Jean gender: F

Row number: [6148] name: Jean gender: F

Row number: [6467] name: Jean gender: F

-----

Row number: [166] name: Dana gender: F

Row number: [1080] name: Dana gender: F

Row number: [1096] name: Dana gender: F

Row number: [1245] name: Dana gender: F

Row number: [2905] name: Dana gender: F

Row number: [6242] name: Dana gender: M

Row number: [6366] name: Dana gender: M

Row number: [7338] name: Dana gender: F

Row number: [7932] name: Dana gender: F

Row number: [8139] name: Dana gender: F

Row number: [8639] name: Dana gender: M

Row number: [8796] name: Dana gender: F

-----

Row number: [228] name: Corey gender: M

Row number: [245] name: Corey gender: F

Row number: [4012] name: Corey gender: F

Row number: [6260] name: Corey gender: F

-----

Row number: [231] name: Dale gender: M

Row number: [2095] name: Dale gender: M

Row number: [2447] name: Dale gender: F

Row number: [6363] name: Dale gender: F

Row number: [7328] name: Dale gender: M

Row number: [7831] name: Dale gender: M

Row number: [8193] name: Dale gender: M

Row number: [9014] name: Dale gender: M

-----

Row number: [301] name: Erin gender: F

Row number: [1



- A = {1,1,2,2,5}; B = {1,2,2,2,5,5}; C = {1,2,3,4,5}

A and B

- J bag sim = 4/11 = 0.3636
- J distance = 7/11 = 0.6363

A and C

- J bag sim = 3/10 = 0.30
- J distance = 7/10 = 0.70

B and C

- J bag sim = 3/11 = 0.2727
- J distance = 8/11 = 0.7272

## Question 9

- a r n a b
- u r b a n

Jaro

- c=2
- t=0
- l=5

$$0.6 = (1/3) * ((2/5)+(2/5))+((2-0)/2))$$

Jaro Winkler

- P=0.1
- L=0

$$0.6 = 0.6 + 0.1 \times 0 \times (0.4)$$

They give the same result due to not having common prefix in both strings. If e.g. they would have a common prefix then Jaro Winkler would give a better matching result.

## Question 10

```
In [208]: sentence = 'Many problems can be expressed as finding similar sets'
k = 5
cardinality_result = len(sentence) + 1 - k
print(f'A set with {k}-shingles has {cardinality_result} cardinality')
```

A set with 5-shingles has 50 cardinality

## Question 11

- D1 = {aa, bb, ab, ba}
- D2 = {aa, ac, ca, ba}
- D3 = {ab, ba, ca}

### a

matrix representation of shingles documents relationship, the universal set is alphabetically ordered

.	D1	D2	D3
aa	1	1	0
ab	1	0	1
ac	0	1	0
ba	1	1	1
bb	1	0	0
ca	0	1	1

### b

index	P1	D1	D2	D3
1	aa	1	1	0
2	bb	1	0	0
3	ab	1	0	1
4	ba	1	1	1
5	ac	0	1	0
6	ca	0	1	1

index	P2	D1	D2	D3
1	ca	0	1	1
2	ac	0	1	0
3	ba	1	1	1
4	ab	1	0	1
5	bb	1	0	0
6	aa	1	1	0

index	P3	D1	D2	D3
1	ac	1	1	0
2	ca	1	0	1
3	ab	0	1	0
4	ba	1	1	1
5	bb	1	0	0
6	aa	0	1	1

signature matrix

.	P1	P2	P3
P1	1	1	3
P2	3	1	1
P3	1	1	2

### C

Comparison of pair documents with their signatures by using Jaccard similarity

.	(D1, D2)	(D1, D3)	(D2, D3)
Col/Col	2/6=0.33	2/5=0.4	2/5=0.4
Sig/Sig	2/3=0.66	0/3=0	1/3=0.33