

## Workshop Week 12 - COMP20008 2020

### Questions

1. Consider the quasi-identifier {job,birth,postcode}.

Job	Birth	Postcode	Illness
Cat1	*	4350	HIV
Cat1	*	4350	HIV
Cat1	1955	5432	flu
Cat1	1955	5432	fever
Cat2	1975	4350	flu
Cat2	1975	4350	fever

- (a) What is the highest  $k$  for which this data is  $k$ -anonymous.
- (b) Describe one possible privacy attack on this data
2. Consider the quasi-identifier {gender,date of birth,zipcode}. Apply generalisation to the following table to make it 3 anonymous.

Name	Gender	Date of birth	ZIP code	Disease
Alice	F	01/01/1981	11111	Flu
Anne	F	02/02/1981	11122	Flu
Sonia	F	12/03/1981	11133	Flu
Bob	M	12/01/1982	33311	Heart disease
Shunsuke	M	10/04/1982	33322	Cold
Carl	M	02/03/1982	33333	Flu

3. Consider the quasi identifier {Age,Zip} for the table below.

Age	Zip	Diagnosis
[21–28]	9****	Measles
[21–28]	9****	Flu
[21–28]	9****	Flu
[48–55]	92***	Cancer
[48–55]	92***	Obesity
[48–55]	92***	Obesity

- (a) What is the highest  $k$  for which this data is  $k$ -anonymous?
  - (b) What is the highest  $l$  for which this data is  $l$ -diverse?
  - (c) Describe one possible privacy attack on this data
4. In the context of providing differential privacy:
    - What is global sensitivity  $G$ ? What is the privacy budget  $k$ ?
    - How does the  $G/k$  ratio affect the noise level?
  5. Consider a survey that collects two values from the respondents, e.g., marital status and sex.
    - Consider a query that takes the survey database as input and outputs a pair of counts (CountNumberFemale, CountNumberMarried). How much can adding or removing an individual affect the output? What is the global sensitivity?
    - Consider a query that takes the survey database as input and outputs the quadruplet of counts (CountMaleMarried, CountMaleSingle, CountFemaleMarried, CountFemaleSingle). How much can adding or removing an individual affect the output? What is the global sensitivity?

## Other Previous Exam Questions

### Exam 2018 - Question 4

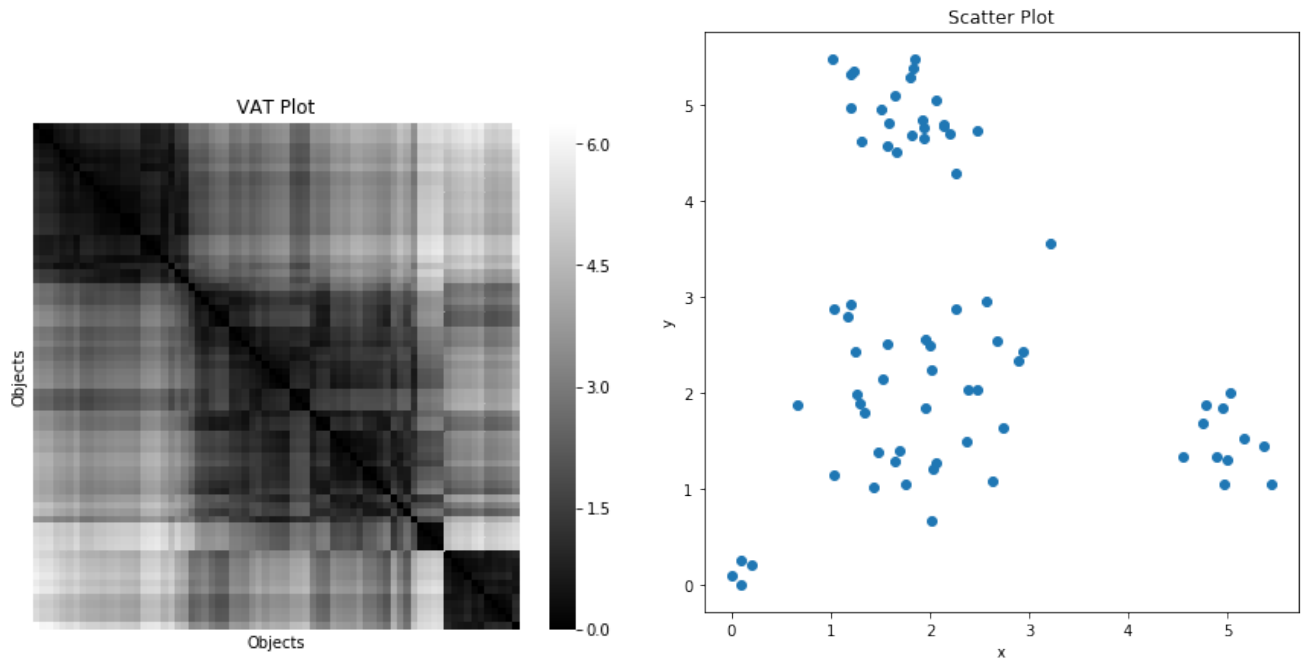
University X is planning to build a recommender system for its students. Based on subjects they have enrolled in, the system will recommend new subjects they might consider studying in the future.

A table showing a fragment of the data input to the system is below. The columns correspond to the codes of all the subjects in the University handbook. Rows correspond to students and whether they have enrolled in a subject (“Yes” if they have previously enrolled and “-” otherwise). The dataset covers the period 2010-2017, with 100,000 students and 3000 subjects. It is proposed that subject recommendations should be made using user based collaborative filtering.

Explain three challenges for making this recommendation approach effective.

PersonName	Subject1	Subject2	Subject3	Subject4	Subject5	...
Alice	Yes	-	-	-	-	...
Bob	Yes	Yes	-	Yes	-	...
Margaret	Yes	Yes	-	-	-	...
...	...	...	...	...	...	...

### Exam 2019 - Question 3



- (a) (1.5 mark) Are these likely to have come from the same dataset? Give two reasons why or why not.
- (b) (1.5 marks) Consider a dataset with 10000 rows and 500 features. Give three reasons why we might want to apply PCA while analysing the dataset.

### Exam 2019 - Question 2b

(3 marks) Consider the following regular expression:

```
int [1-9] += [A-Z] [a-z] *
```

Which of the following terms would be matched by this regular expression (write down all that apply):

```
int=Abc
int2=A
int3=Abc
int44=Azz
int5Abc
int6=abc
```