# Attacks on k-anonymous table

| Table 1.2. Generalized medical record table. |    |         |        |          |             |                 |
|--|----|---------|--------|----------|-------------|-----------------|
|  |    | Age     | Gender | Zip Code | Nationality | Condition       |
| (Ann)  | 1  | 20-29   | Any    | 130**    | Any         | Heart disease   |
| (Bruce)                                      | 2  | 20 - 29 | Any    | 130**    | Any         | Heart disease   |
| (Cary)                                       | 3  | 20 - 29 | Any    | 130**    | Any         | Viral infection |
| (Dick)                                       | 4  | 20 - 29 | Any    | 130**    | Any         | Viral Infection |
| (Eshwar)                                     | 5  | 40 - 59 | Any    | 14***    | Asian       | Cancer          |
| (Fox)  | 6  | 40 - 59 | Any    | 14***    | Asian       | Flu             |
| (Gary)                                       | 7  | 40 - 59 | Any    | 14***    | Asian       | Heart disease   |
| (Helen)                                      | 8  | 40 - 59 | Any    | 14***    | Asian       | Flu             |
| (Igor)                                       | 9  | 30-39   | Any    | 1322*    | American    | Cancer          |
| (Jean)                                       | 10 | 30 - 39 | Any    | 1322*    | American    | Cancer          |
| (Ken)  | 11 | 30-39   | Any    | 1322*    | American    | Cancer          |
| (Lewis)                                      | 12 | 30-39   | Any    | 1322*    | American    | Cancer          |

## Attacks on k-anonymous table

| Table 1.2. Generalized medical record table. |    |         |        |          |             |                 |
|--|----|---------|--------|----------|-------------|-----------------|
| -  |    | Age     | Gender | Zip Code | Nationality | Condition       |
| (Ann)  | 1  | 20-29   | Any    | 130**    | Any         | Asthma          |
| (Bruce)                                      | 2  | 20 - 29 | Any    | 130**    | Any         | Asthma          |
| (Cary)                                       | 3  | 20 - 29 | Any    | 130**    | Any         | Viral infection |
| (Dick)                                       | 4  | 20 - 29 | Any    | 130**    | Any         | Viral Infection |
| (Eshwar)                                     | 5  | 40 - 59 | Any    | 14***    | Asian       | Cancer          |
| (Fox)  | 6  | 40 - 59 | Any    | 14***    | Asian       | Flu             |
| (Gary)                                       | 7  | 40 - 59 | Any    | 14***    | Asian       | Heart disease   |
| (Helen)                                      | 8  | 40 - 59 | Any    | 14***    | Asian       | Flu             |
| (Igor)                                       | 9  | 30-39   | Any    | 1322*    | American    | Cancer          |
| (Jean)                                       | 10 | 30 - 39 | Any    | 1322*    | American    | Cancer          |
| (Ken)  | 11 | 30 - 39 | Any    | 1322*    | American    | Cancer          |
| (Lewis)                                      | 12 | 30-39   | Any    | 1322*    | American    | Cancer          |

#### Prior and posterior beliefs

| Table 1.2. | Generalized | medical | record | table. |
|------------|-------------|---------|--------|--------|
|            |             |         |        |        |

| s)       |    | Age     | Gender | Zip Code | Nationality | Condition       |
|----------|----|---------|--------|----------|-------------|-----------------|
| (Ann)    | 1  | 20-29   | Any    | 130**    | Any         | Heart disease   |
| (Bruce)  | 2  | 20 - 29 | Any    | 130**    | Any         | Heart disease   |
| (Cary)   | 3  | 20 - 29 | Any    | 130**    | Any         | Viral infection |
| (Dick)   | 4  | 20 - 29 | Any    | 130**    | Any         | Viral Infection |
| (Eshwar) | 5  | 40-59   | Any    | 14***    | Asian       | Cancer          |
| (Fox)    | 6  | 40 - 59 | Any    | 14***    | Asian       | Flu             |
| (Gary)   | 7  | 40 - 59 | Any    | 14***    | Asian       | Heart disease   |
| (Helen)  | 8  | 40 – 59 | Any    | 14***    | Asian       | Flu             |
| (Igor)   | 9  | 30-39   | Any    | 1322*    | American    | Cancer          |
| (Jean)   | 10 | 30 - 39 | Any    | 1322*    | American    | Cancer          |
| (Ken)    | 11 | 30 - 39 | Any    | 1322*    | American    | Cancer          |
| (Lewis)  | 12 | 30 - 39 | Any    | 1322*    | American    | Cancer          |

Prior knowledge  $\alpha = p(t[S] = s \mid t[Q] = q)$ 

Posterior knowledge  $\beta = p(t[S] = s | t[Q] = q$   $\land \exists t^* \in T^*, t \rightarrow t^*)$ 

## How to prevent homogeneity attacks?

|    | l        | Von-Sen   | Sensitive   |                 |
|----|----------|-----------|-------------|-----------------|
|    | Zip Code | Age       | Nationality | Condition       |
| 1  | 130**    | < 30      | *           | Heart Disease   |
| 2  | 130**    | < 30      | *           | Heart Disease   |
| 3  | 130**    | < 30      | *           | Viral Infection |
| 4  | 130**    | < 30      | *           | Viral Infection |
| 5  | 1485*    | $\geq 40$ | *           | Cancer          |
| 6  | 1485*    | $\geq 40$ | *           | Heart Disease   |
| 7  | 1485*    | $\geq 40$ | *           | Viral Infection |
| 8  | 1485*    | $\geq 40$ | *           | Viral Infection |
| 9  | 130**    | 3*        | *           | Cancer          |
| 10 | 130**    | 3*        | *           | Cancer          |
| 11 | 130**    | 3*        | *           | Cancer          |
| 12 | 130**    | 3*        | *           | Cancer          |

#### ℓ-diversity

An equivalence class (E) is  $\ell$ -diverse if it contains at least  $\ell$  well-represented values for the sensitive attribute S.

A table is  $\ell$ -diverse if every E is  $\ell$ -diverse.

|    | 1        | Von-Sen   | Sensitive   |                 |
|----|----------|-----------|-------------|-----------------|
|    | Zip Code | Age       | Nationality | Condition       |
| 1  | 130**    | < 30      | *           | Heart Disease   |
| 2  | 130**    | < 30      | *           | Heart Disease   |
| 3  | 130**    | < 30      | *           | Viral Infection |
| 4  | 130**    | < 30      | *           | Viral Infection |
| 5  | 1485*    | $\geq 40$ | *           | Cancer          |
| 6  | 1485*    | $\geq 40$ | *           | Heart Disease   |
| 7  | 1485*    | $\geq 40$ | *           | Viral Infection |
| 8  | 1485*    | $\geq 40$ | *           | Viral Infection |
| 9  | 130**    | 3*        | *           | Cancer          |
| 10 | 130**    | 3*        | *           | Cancer          |
| 11 | 130**    | 3*        | *           | Cancer          |
| 12 | 130**    | 3*        | *           | Cancer          |

#### Distinct ℓ-diversity

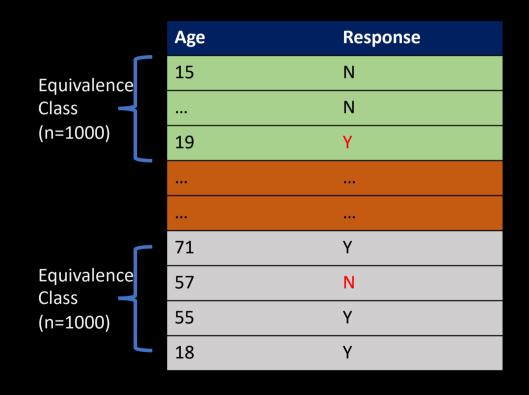
At least & distinct values for the sensitive attribute in each equivalence class.

|    | l l      | Von-Sen   | Sensitive   |                 |
|----|----------|-----------|-------------|-----------------|
|    | Zip Code | Age       | Nationality | Condition       |
| 1  | 130**    | < 30      | *           | Heart Disease   |
| 2  | 130**    | < 30      | *           | Heart Disease   |
| 3  | 130**    | < 30      | *           | Viral Infection |
| 4  | 130**    | < 30      | *           | Viral Infection |
| 5  | 1485*    | $\geq 40$ | *           | Cancer          |
| 6  | 1485*    | $\geq 40$ | *           | Heart Disease   |
| 7  | 1485*    | $\geq 40$ | *           | Viral Infection |
| 8  | 1485*    | $\geq 40$ | *           | Viral Infection |
| 9  | 130**    | 3*        | *           | Cancer          |
| 10 | 130**    | 3*        | *           | Cancer          |
| 11 | 130**    | 3*        | *           | Cancer          |
| 12 | 130**    | 3*        | *           | Cancer          |

|    | N        | Von-Sen   | Sensitive   |                 |
|----|----------|-----------|-------------|-----------------|
|    | Zip Code | Age       | Nationality | Condition       |
| 1  | 1305*    | $\leq 40$ | *           | Heart Disease   |
| 4  | 1305*    | $\leq 40$ | *           | Viral Infection |
| 9  | 1305*    | $\leq 40$ | *           | Cancer          |
| 10 | 1305*    | $\leq 40$ | *           | Cancer          |
| 5  | 1485*    | > 40      | *           | Cancer          |
| 6  | 1485*    | > 40      | *           | Heart Disease   |
| 7  | 1485*    | > 40      | *           | Viral Infection |
| 8  | 1485*    | > 40      | *           | Viral Infection |
| 2  | 1306*    | $\leq 40$ | *           | Heart Disease   |
| 3  | 1306*    | $\leq 40$ | *           | Viral Infection |
| 11 | 1306*    | $\leq 40$ | *           | Cancer          |
| 12 | 1306*    | $\leq 40$ | *           | Cancer          |

### Skewness within equivalence classes

| Age | Response<br>(sensitive) |
|-----|-------------------------|
| 15  | Υ                       |
| 21  | N                       |
| 19  | Υ                       |
| 66  | N                       |
| 58  | Υ                       |
| 71  | N                       |
|     |                         |
| 55  | N                       |



100,000 records

#### Entropy \(\ell\)-diversity

The entropy of an equivalent class E is defined to be

$$H(E) = -\sum_{s \in S} p(E, s) \log(p(E, s))$$

An equivalence class E has entropy ℓ-diversity if

$$H(E) \ge \log(l)$$

A table is  $\ell$ -diverse if each E is  $\ell$ -diverse.

#### Entropy \(\ell\)-diversity

The entropy of an equivalent class E is defined to be

$$H(E) = -\sum_{s \in S} p(E, s) \log(p(E, s))$$

An equivalence class E has entropy ℓ-diversity if

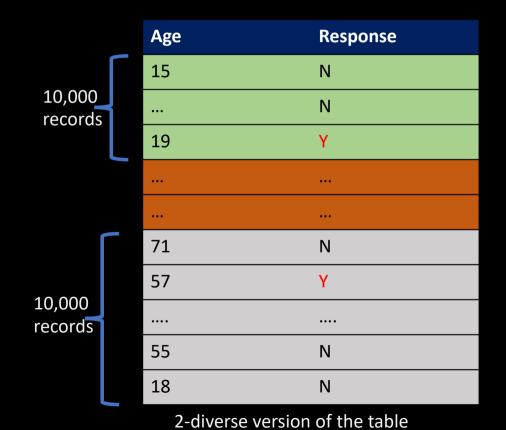
$$H(E) \ge \log(l)$$

Constraint: the full table must have  $H(T) \ge \log(l)$  because  $H(T) \ge \min(H(E_a), H(E_b))$ 

#### Skewness in the population

| Age | Response |
|-----|----------|
| 15  | N        |
| 21  | N        |
| 19  | Υ        |
| 66  | N        |
| 58  | N        |
| 71  | N        |
|     |          |
| 55  | N        |

100,000 Records, 99% negative



#### Recursive $(c,\ell)$ -diversity

Let m be the number of values in an equivalence class, and  $r_i$   $(1 \le i \le m)$  be the number of times that the *ith* most frequent sensitive value appears in an equivalence class **E**.

Then, E has recursive (c,  $\ell$ )-diversity if  $r_1 < c(r_l, r_{l+1} + ... + r_m)$ 

A table is said to have recursive  $(c, \ell)$ -diversity if all of its equivalence classes have recursive  $(c, \ell)$ -diversity.

#### Recursive $(c,\ell)$ -diversity

| Table 1.2. | Generalized   | medical  | record   | table. |
|------------|---------------|----------|----------|--------|
| 10010 1.2. | GCIICI GIIZCG | mountain | 1 CCOI G | UUUDIC |

|          |    | Age     | Gender | Zip Code | Nationality | Condition       |
|----------|----|---------|--------|----------|-------------|-----------------|
| (Ann)    | 1  | 20 - 29 | Any    | 130**    | Any         | Heart disease   |
| (Bruce)  | 2  | 20 - 29 | Any    | 130**    | Any         | Heart disease   |
| (Cary)   | 3  | 20-29   | Any    | 130**    | Any         | Viral infection |
| (Dick)   | 4  | 20 - 29 | Any    | 130**    | Any         | Viral Infection |
| (Eshwar) | 5  | 40-59   | Any    | 14***    | Asian       | Cancer          |
| (Fox)    | 6  | 40 – 59 | Any    | 14***    | Asian       | Flu             |
| (Gary)   | 7  | 40 – 59 | Any    | 14***    | Asian       | Heart disease   |
| (Helen)  | 8  | 40–59   | Any    | 14***    | Asian       | Flu             |
| (Igor)   | 9  | 30-39   | Any    | 1322*    | American    | Cancer          |
| (Jean)   | 10 | 30 - 39 | Any    | 1322*    | American    | Cancer          |
| (Ken)    | 11 | 30-39   | Any    | 1322*    | American    | Cancer          |
| (Lewis)  | 12 | 30 - 39 | Any    | 1322*    | American    | Cancer          |

Then, E has recursive  $(c, \ell)$ -diversity if  $r_1 < c(r_2 + r_3)$ 

#### Recursive $(c,\ell)$ -diversity

$$r_1 < c(r_l, r_{l+1} + \dots + r_m)$$

The count of the most frequent item must be *less than* a constant multiple of the total count of the least (m-l) frequent items.

#### Limitations of $\ell$ -diversity: drastic loss of utility

| Zip   | Disease |
|-------|---------|
| 15152 | Υ       |
| 15163 | Υ       |
| 15784 | Υ       |
| 18784 | Ν       |
| 18654 | N       |
| 18744 | N       |

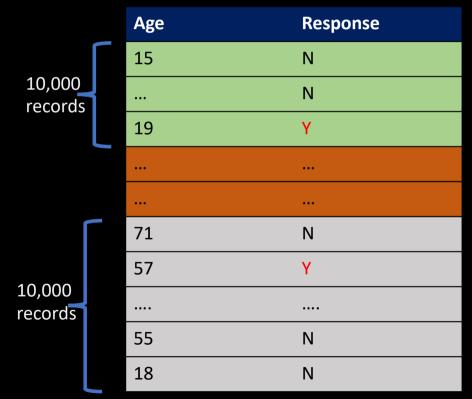
| Zip   | Disease |
|-------|---------|
| 15152 | Υ       |
| 18784 | N       |
| 15163 | Υ       |
| 18784 | N       |
| 15784 | Υ       |
| 18744 | N       |

| zip   | Disease |
|-------|---------|
| 1**** | Υ       |
| 1**** | N       |
| 1**** | Υ       |
| 1**** | N       |
| 1**** | Υ       |
| 1**** | N       |

### Sometimes \(\ell\)-diversity may be unnecessary

| Age | Response |
|-----|----------|
| 15  | N        |
| 21  | N        |
| 19  | Υ        |
| 66  | N        |
| 58  | N        |
| 71  | N        |
|     |          |
| 55  | N        |

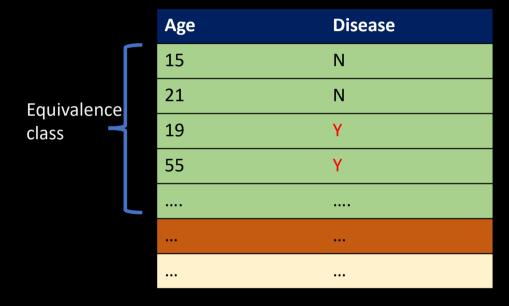
100,000 Records, 99% negative



2-diverse version of the table

# Limitations of $\ell$ -diversity: wrong inference because of skewness

| Age | Disease |
|-----|---------|
| 15  | N       |
| 21  | N       |
| 19  | Υ       |
| 66  | N       |
| 58  | N       |
| 71  | N       |
|     |         |
| 55  | Υ       |



2-diverse version of the table

100,000 Records, 99% negative

#### Limitations of $\ell$ -diversity: similarity attack

| Age | Disease        |
|-----|----------------|
| 15  | Lung cancer    |
| 21  | Stomach cancer |
| 19  | No issue       |
| 66  | Stomach cancer |
| 58  | Liver cancer   |
| 71  | No issue       |
|     |                |
| 55  | Headache       |

| Age | Response       |
|-----|----------------|
| 15  | Lung cancer    |
|     | Stomach cancer |
| 19  | Liver cancer   |
| :   |                |
| :   |                |
| 71  |                |
|     |                |
| 55  |                |
| 18  |                |

#### Intuition behind the limitations of I-diversity

Intuitively, distributions that have the same level of (syntactic) diversity may provide very different levels of privacy, because

- 1. there are semantic relationships among the attribute value
- 2. different values have very different levels of sensitivity, e.g., headache can have vastly different level of sensitivity than cancer/diabetic
- 3. different distributions at the equivalent class level and population level can have adverse effects

#### Intuition behind t-closeness

| Table 1.2. Generalized medical record table. |    |         |        |          |             |                 |
|--|----|---------|--------|----------|-------------|-----------------|
|  |    | Age     | Gender | Zip Code | Nationality | Condition       |
| (Ann)  | 1  | 20-29   | Any    | 130**    | Any         | Heart disease   |
| (Bruce)                                      | 2  | 20 - 29 | Any    | 130**    | Any         | Heart disease   |
| (Cary)                                       | 3  | 20 - 29 | Any    | 130**    | Any         | Viral infection |
| (Dick)                                       | 4  | 20 - 29 | Any    | 130**    | Any         | Viral Infection |
| (Eshwar)                                     | 5  | 40 - 59 | Any    | 14***    | Asian       | Cancer          |
| (Fox)  | 6  | 40 - 59 | Any    | 14***    | Asian       | Flu             |
| (Gary)                                       | 7  | 40 - 59 | Any    | 14***    | Asian       | Heart disease   |
| (Helen)                                      | 8  | 40 - 59 | Any    | 14***    | Asian       | Flu             |
| (Igor)                                       | 9  | 30-39   | Any    | 1322*    | American    | Cancer          |
| (Jean)                                       | 10 | 30 - 39 | Any    | 1322*    | American    | Cancer          |
| (Ken)  | 11 | 30 - 39 | Any    | 1322*    | American    | Cancer          |
| (Lewis)                                      | 12 | 30–39   | Any    | 1322*    | American    | Cancer          |

| Table 1.2. Generalized medical record table.              |     |        |          |             |  |
|---|-----|--------|----------|-------------|--|
|   | Age | Gender | Zip Code | Nationality | Condition  |
| (Ann) (Bruce) (Cary) (Dick) (Eshwar) (Fox) (Gary) (Helen) |     |        |          | Q $=$       | Heart disease Heart disease Viral infection Viral Infection Cancer Flu Heart disease Flu |
| (Igor) (Jean) (Ken) (Lewis)                               |     | v      |          |             | Cancer<br>Cancer<br>Cancer<br>Cancer   |

#### Intuition behind t-closeness

| Table 1.2. Generalized medical record table. |    |         |        |            |             |                 |
|--|----|---------|--------|------------|-------------|-----------------|
|  |    | Age     | Gender | Zip Code   | Nationality | Condition       |
| (Ann)  | 1  | 20-29   | Any    | 130**      | Any         | Heart disease   |
| (Bruce)                                      | 2  | 20 - 29 | Any    | 130**      | Any         | Heart disease   |
| (Cary)                                       | 3  | 20 - 29 | Any    | 130**      | Any         | Viral infection |
| (Dick)                                       | 4  | 20 - 29 | Any    | 130**      | Any         | Viral Infection |
| (Eshwar)                                     | 5  | 40 – 59 | Any    | 14***      | Asian       | Cancer          |
| (Fox)  | 6  | 40 - 59 | Any    | 14***      | Asian       | Flu             |
| (Gary)                                       | 7  | 40 - 59 | Any    | $14^{***}$ | Asian       | Heart disease   |
| (Helen)                                      | 8  | 40 – 59 | Any    | 14***      | Asian       | Flu             |
| (Igor)                                       | 9  | 30-39   | Any    | 1322*      | American    | Cancer          |
| (Jean)                                       | 10 | 30 - 39 | Any    | 1322*      | American    | Cancer          |
| (Ken)  | 11 | 30 - 39 | Any    | 1322*      | American    | Cancer          |
| (Lewis)                                      | 12 | 30-39   | Any    | 1322*      | American    | Cancer          |
|  |    |         |        |            |             |                 |

Prior belief about an individual,  $\alpha$ 

Changed belief after seeing distribution at the population level,  $\beta$ 

Changed belief after looking at the (anonymized) table,  $\gamma$ 

l —diversity minimizes  $\gamma - \alpha$ 

t —closeness minimizes  $\gamma - \beta$ 

#### t —closeness minimizes $\gamma - \beta$

| _                                      | Table 1.2. Generalized medical record table. |        |          |             | Q  |
|--|--|--------|----------|-------------|--|
|  | Age  | Gender | Zip Code | Nationality | Condition  |
| (Ann) (Bruce) (Cary) (Dick)            |  |        |          |             | Heart disease<br>Heart disease<br>Viral infection<br>Viral Infection |
| (Eshwar)<br>(Fox)<br>(Gary)<br>(Helen) |  |        |          |             | Cancer<br>Flu<br>Heart disease<br>Flu                                |
| (Igor)<br>(Jean)<br>(Ken)<br>(Lewis)   |  |        |          | _           | Cancer<br>Cancer<br>Cancer<br>Cancer                                 |

| Table 1.2. Generalized medical record table. |      |   |         |        |          |             |                 |
|--|------|---|---------|--------|----------|-------------|-----------------|
|  |      |   | Age     | Gender | Zip Code | Nationality | Condition       |
| (Anı   | n)   | 1 | 20-29   | Any    | 130**    | Any         | Heart disease   |
| (Bruc  | e)   | 2 | 20 - 29 | Any    | 130**    | Any         | Heart disease   |
| (Car   | y)   | 3 | 20 - 29 | Any    | 130**    | Any         | Viral infection |
| (Dicl  | x)   | 4 | 20 – 29 | Any    | 130**    | Any         | Viral Infection |
| (Eshwa                                       | r)   | 5 | 40-59   | Any    | 14***    | Asian       | Cancer          |
| D (For                                       | x)   | 6 | 40 - 59 | Any    | 14***    | Asian       | Flu             |
| (Gar   | y)   | 7 | 40 - 59 | Any    | 14***    | Asian       | Heart disease   |
| (Helei                                       | n)   | 8 | 40 - 59 | Any    | 14***    | Asian       | Flu             |
| (Igo   | r)   | 9 | 30-39   | Any    | 1322*    | American    | Cancer          |
| (Jean  | n) 1 | 0 | 30 - 39 | Any    | 1322*    | American    | Cancer          |
| (Kei   | n) 1 | 1 | 30 - 39 | Any    | 1322*    | American    | Cancer          |
| (Lewi  | s) 1 | 2 | 30–39   | Any    | 1322*    | American    | Cancer          |

Q(heart disease) = 1/4Q(Flu) = 1/6 P(heart disease) = 1/4P(Flu) = 1/2

#### *t*-closeness

An equivalence class has t-closeness if the distance between the distribution of a sensitive attribute in this class and the distribution of the attribute in the whole table is no more than a threshold t.

$$dist(P,Q) \le t$$

A table is said to have t-closeness if all equivalence classes have t-closeness.

#### *t*-closeness

$$dist(P,Q) \leq t$$

#### Distance measures

- 1. Kulback-Liebler divergence
- 2. Jensen-Shannon
- 3. Earth mover's distance

#### resources

**Diversity: Privacy Beyond k-Anonymity** 

t-closeness: Privacy beyond k-anonymity and l-diversity