# Scalable Data Engineering, Exercise 4

Henry Haustein

## Task 1

(a) True

(b) False, they can't be proven algorithmically.

(c) True

## Task 2

The bigrams are

- firstname: f, fi, ir, rs, st, tn, na, am, me, e
- lastname: l, la, as, st, tn, na, am, me, e
- street: s, st, tr, re, ee, et, t
- name: n, na, am, me, e
- address: a, ad, dd, dr, re, es, ss, s
- forename: f, fo, or, re, en, na, am, me, e

Then the similarities are

|          | firstname | lastname | street |
|----------|-----------|----------|--------|
| **name**    | $\frac{2\cdot4}{10+5}=0.533$ | $\frac{2\cdot2\cdot4}{9+5}=0.571$ | $\frac{2\cdot0}{7+5}=0$ |
| **address** | $\frac{2\cdot0}{10+8}=0$ | $\frac{2\cdot0}{9+8}=0$ | $\frac{2\cdot1}{7+8}=0.133$ |
| **forename**| $\frac{2\cdot5}{10+9}=0.526$ | $\frac{2\cdot4}{9+9}=0.444$ | $\frac{2\cdot1}{7+9}=0.125$ |

Stable Marriage Algorithm:

- firstname proposes to name, agrees: (firstname, name)
- lastname proposes to name, agrees + leaves: (lastname, name)
- firstname proposes to forename, agrees: (lastname, name), (firstname, forename)
- street proposes to adress, agrees: (lastname, name), (firstname, forename), (street, adress)

## Task 3

(a)

(b) Nation ⇔ Country, (Region ⇔ Continent)

(c) SQL:

```
1  ALTER TABLE mondial.country ADD COLUMN cid UUID
2
3  UPDATE mondial.country SET cid = gen_random_uuid()
```

(d) SQL:

```
1  ALTER TABLE supplier ADD COLUMN s_countrykey UUID
2
3  UPDATE supplier SET s_countrykey = c.uuid
4  FROM nation AS n, mondial.country AS c
5  WHERE
6    lower(n.name) = lower(c.name) AND
7    supplier.s_nationkey = nation.n_nationkey
```