



Harmonisation methods

Eve Cheng

The goal

A bunch of survey questions:

1. How old are you?
2. What is your occupation?
3. How old is your dog?
4. Who would you vote for?
5. How old is your father?
6. How old is your mother?

?

1. How old are you?
2. How old is your dog?
3. How old is your father?
4. How old is your mother?

1. What is your occupation?

1. Who would you vote for?





Current methods

Machine learning: requires training your model or downloading existing models

word2vec
(word
mover)

SentenceTransformer

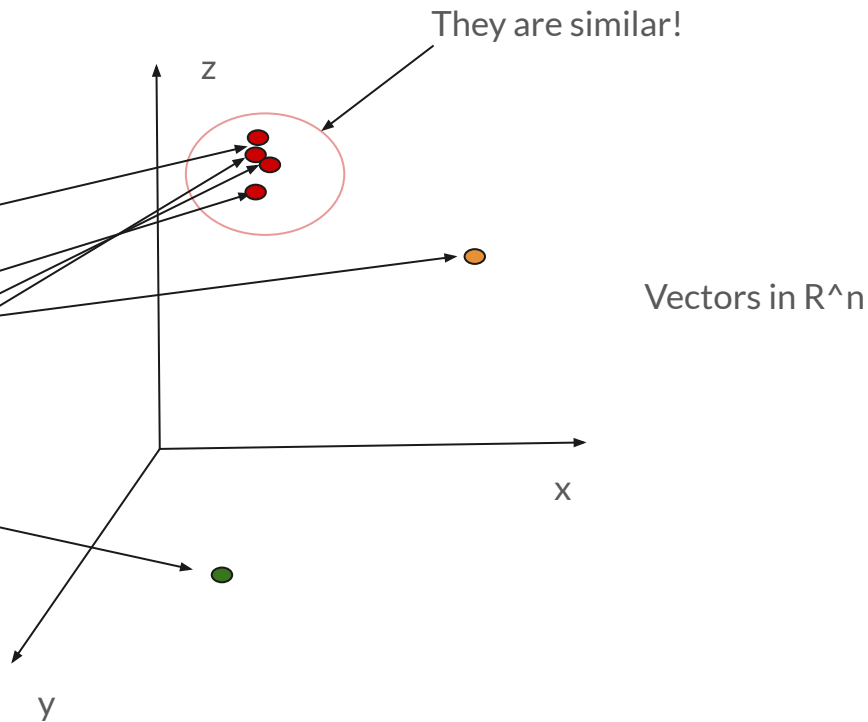
No training required

Bag of words

Normalised BOW

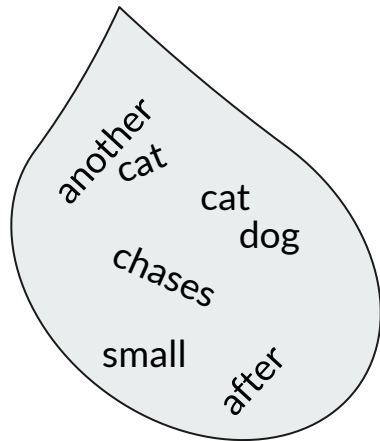
A bunch of survey questions:

1. How old are you?
2. What is your occupation?
3. How old is your dog?
4. Who would you vote for?
5. How old is your father?
6. How old is your mother?



Step 1: calculate the vectors

The cat chases after the dog
and another small cat



word	freq
cat	2
dog	1
chases	1
small	1
after	1
another	1

Vector in \mathbb{R}^6
(2,1,1,1,1,1)

Step 2: calculate pair-wise distances

Sentence 1 → Vector 1

Sentence 2 → Vector 2

Sentence 3 → Vector 3

⋮

Sentence n → Vector n

Similarity matrix (n x n)

$d(v_1, v_1)$	$d(v_1, v_2)$	$d(v_1, v_3)$
$d(v_2, v_1)$	$d(v_2, v_2)$	$d(v_2, v_3)$
$d(v_3, v_1)$	$d(v_3, v_2)$	$d(v_3, v_3)$

...

⋮

$d(v_1, v_0)$: **cosine similarity** between v_1 and v_0

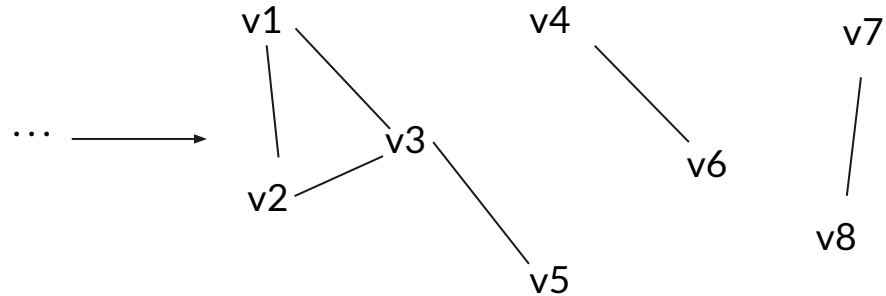
Step 3: create a network of sentences

Similarity matrix (n x n)

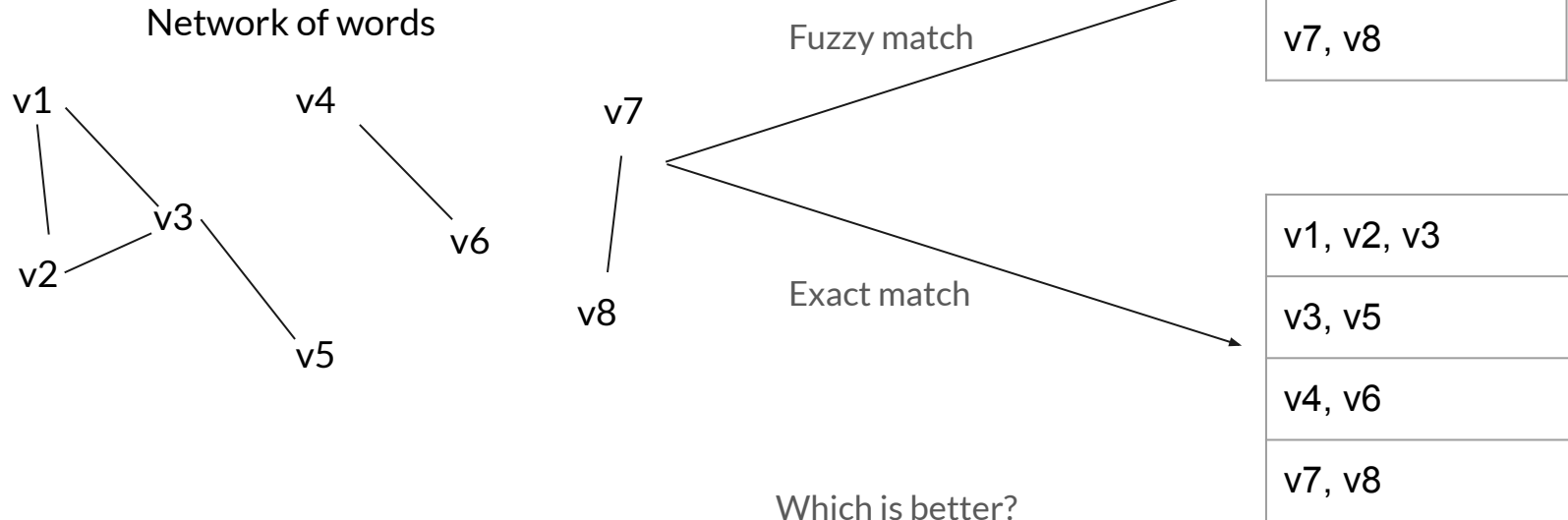
$d(v1, v1)$	$d(v1, v2)$	$d(v1, v3)$
$d(v2, v1)$	$d(v2, v2)$	$d(v2, v3)$
$d(v3, v1)$	$d(v3, v2)$	$d(v3, v3)$

⋮

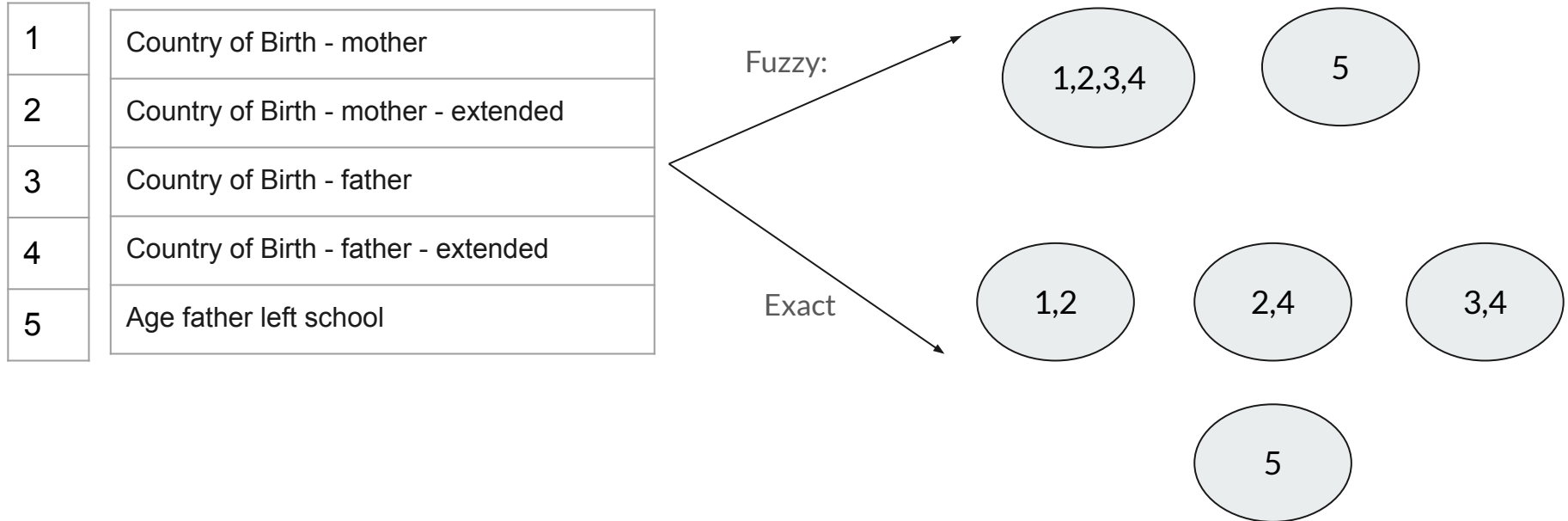
Network of words



Step 4: detect clusters



Example. Input 2: threshold (set to 0.7 here)





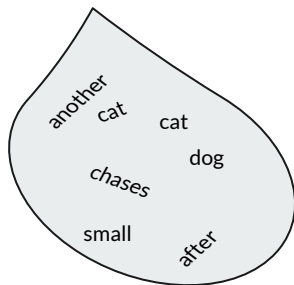
Inputs for BOW so far:

1. Threshold: above what similarity do you consider similar. The program requires a number from 0 to 1.
2. Cluster option: fuzzy or exact

One last input: the N in the Ngram method

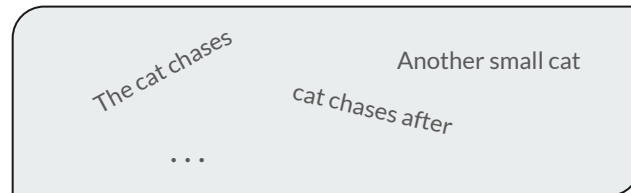
Vanilla BOW

The cat chases after the
dog and another small cat



Ngram: for example, N=3

The cat chases after the dog and another small cat.





Summary

Threshold

Cluster option

N in Ngram



```
graph LR; Threshold --> BOWpy[BOW.py]; Cluster_option[Cluster option] --> BOWpy; N_in_Ngram[N in Ngram] --> BOWpy; BOWpy --> Output[Sentences that are deemed similar];
```

BOW.py

Sentences that are deemed similar



What about all those machine learning stuff?

