

Metanome

For this question I used the Metanome tool together with the FastFDs-1.2 and DepMiner-1.2 algorithm.

For FastFD (19 results), please see fd_algorithm_1.png and fd_algorithm_2.png

determinant	dependant
Abbreviation	Name
Unicode	Name
UnicodeDisplay	Name
Represents	Name
Name	AU_Abbreviation
Unicode	AU_Abbreviation
UnicodeDisplay	AU_Abbreviation
Represents	AU_Abbreviation
Name	Symbol
AU_Abbreviation	Symbol
Unicode	Symbol
UnicodeDisplay	Symbol
Represents	Symbol
UnicodeDisplay	Unicode
Represents	Unicode
Unicode	UnicodeDisplay
Represents	UnicodeDisplay
Unicode	Represents
UnicodeDisplay	Represents

For DepMiner (14 results), please see dm_algorithm_1.png and dm_algorithm_2.png

determinant	dependant
IAU_Abbreviation	Name
UnicodeDisplay	Name
Represents	Name
Unicode	Name
Represents	UnicodeDisplay
Unicode	UnicodeDisplay
UnicodeDisplay	IAU_Abbreviation
Represents	IAU_Abbreviation
Name	IAU_Abbreviation
Unicode	IAU_Abbreviation
UnicodeDisplay	Unicode
Represents	Unicode
UnicodeDisplay	Represents
Unicode	Represents

Algorithm FastFD (19 results) has more results than DepMiner (14 results).

Question 6

Computing Levenshtein distance

"*" = empty space

"*" = shows the minimum cost path

1

Computed

	.	C	o	m	p	l	e	t	i	o	n
.	+	0	1	2	3	4	5	6	7	8	9
C	1	+	0	1	2	3	4	5	6	7	8
o	2	1	+	0	1	2	3	4	5	6	7
m	3	2	1	+	0	1	2	3	4	5	6
p	4	3	2	1	+	0	1	2	3	4	5
u	5	4	3	2	1	+	0	1	2	3	4
t	6	5	4	3	2	3	+	4	3	4	5
a	7	6	5	4	3	4	4	+	5	4	5
t	8	7	6	5	4	5	5	4	+	5	6
i	9	8	7	6	5	6	6	5	4	+	5
o	10	9	8	7	6	7	7	6	5	4	+
n	11	10	9	8	7	8	8	7	6	5	4

be * between p and l of completion)

= 5 cost

ce to use less cost)

ost out of the edit distances (insert, substitute, delete)
e cost for e.g. insert/delete in the edit distances at

nce

result = 3 cost

Handwritten

- Computation
- Comp*letion

d+s+s = 3 operations = 3 cost

2

Computed

Levenshtein-Winkler similarity

result = 5 cost

Handwritten

- Computation
- Comp*letion

d+s+s = 3 operations = i + (2 x 2) = 5 cost

3

Computed

		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	

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
-------	----	----	----	----

result = 5 cost

Handwritten

- Computation
- Comp*letion

d+s+s = 3 operations = i + (2 x 3) = 7 cost

(this is the result if we keep the empty space * between p and l of completion)

- Comp**utation
- Comp***tion

i+i+d+d+d = 5 operations = (i x 1) + (d x 1) = 5 cost

(this is the result if we add more empty space to use less cost)

When the update cost is > 3:

The algorithm is looking for the minimum cost out of the edit distances (insert, substitute, deletion) from the current position [i,j]. So, if the cost for substitution is high but the cost for e.g. insert/ delete in the edit distances at [i,j] position is lower, than insert/ delete operation is applied.

Question 7

Computing the gap distance

- (i) insertion cost = 1
- (o) open gap cost = 1
- (e) extend gap cost = 0.1

Adv ances in Instrum entation and Control

Adv. _____ Instrum. _____ Control

Adv[i][o][e][e][e][e][e][e][e] Instrum[i][o][e][e][e][e][e][e][e] Control

- i = 1 + 1 = 2
- o = 1 + 1 = 2
- e = 8 + 12 = 20

cost = i+o+e = 2 + 2 + (20 x 0.1)

cost = 24

Question 8

Compute Jaccard bag similarity and Jaccard Distance

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

Compute the Jaro and Jaro-Winkler similarity

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

Jaro sim

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

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

Jaro Winkler

- P=0.1
- L=0

Jw(S1, S2) = JaroSim + P x L x (1-JaroSim)

0.6 = 0.6 + 0.1 x 0 x (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

Cardinality of the set that has 5-shingles

```
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

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 - signature matrix

Permutation 1

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

Permutation 2

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

Permutation 3

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	1	1	3
P2	3	1	1
P3	1	1	2

c - Jaccard similarity

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

Info on col and sig:

- D1 and D2: on 'Col' similarity is 33% and on 'Sig' the agreement is 66%.
- D1 and D3: on 'Col' similarity is 40% and on 'Sig' the agreement is 0%.
- D2 and D3: on 'Col' similarity is 40% and on 'Sig' the agreement is 33%.