

Approximate
String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Approximate String Matching

Assignment Project Exam Help

<https://eduassistpro.github.io>

Jeremy Nicholson and Justin Zobel and Karin Verspoor

Add WeChat edu_assist_pr

Search



THE UNIVERSITY OF

MELBOURNE

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Week 3:

- Approximate String Search and Matching



- Edit Distance

- N-Gram Distance

- [Phonetic methods]



- Evaluation



- [Genomics]

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Consider:

<https://eduassistpro.github.io>

In exes for foxes rex dux mixes a pox of waxed luxe

An axe, and an axon, to exo Exxon max oxen.

Grexit or Brexit as quixotic bankers with b

Add WeChat edu_assist_pro

Not (really) a Knowledge Technology!

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Find exon in:

<https://eduassistpro.github.io>

Not present!

...But what is the “closest” or “best” match?

This is a Knowledge Technology!

Add WeChat edu_assist_pro

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

Need the notion of a **dictionary**:



misspelled



An item in the input which *does*
correctly spelled or misspelled
(of this subject)

Add WeChat edu_assist_pr

<https://eduassistpro.github.io>

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

Depends on the person who wrote the origin

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

- Computational Genomics (later, if we have time)

■
■
<https://eduassistpro.github.io>

- Data cleaning

■
Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

Grexit or Brexit as quixotic haxxers with b

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

For a given string w of interest:



s

<https://eduassistpro.github.io>



All results found in dictionary are returned

Unix command-line utility `fgrep` is an
these. Add WeChat `edu_assist_pr`

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

Requires 1 insertion (o) so intended wo
neighbourhood search (and some uninte

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

Neighbourhood search is surprisingly fast!

<https://eduassistpro.github.io>

...But Σ is a small constant, string of interest is usually short, and k is usually small

Add WeChat [edu_assist_pro](#)

For each neighbour, need a dictionary read
Binary search yields $\mathcal{O}(|w|^k \log D)$

Approximate
String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

Each operation is associated with a score;

Best match is the dictionary entry with best ag

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

For example:

Item of interest: `crat`

<https://eduassistpro.github.io>

`crat` → `cart`:

Match `c` (+1), Delete `r` (-1), Match `a`

+1) = +1

`crat` → `arts`:

Replace `c` with `a` (-1), Match `r` (+1),

(-1) = -1

sert `s`

`cart` is the better match

Add WeChat `edu_assist_pro`

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

This is the Levenshtein Distance (which is a “
number of edits required to transform one str

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

Hypothetically, any parameter is possible!

<https://eduassistpro.github.io>

aba

- foo: Insert, Delete, Insert, Delete, Ins

- aba: Match, Match, Match = +12

- cbc: Replace, Match, Replace = +4

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

Consider:

Is faxing more likely to be facing

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

From string f to string t , given array A of $|f| + 1$ columns and $|t| + 1$ rows, we can solve using the Needleman–Wunsch algorithm:

```
lf = strlen(f); lt = strlen(t);  
A[0][0]=0;
```

```
for (k=1; k<=lf; k++)
```

```
    A[j][k] = max3( //Or min3 if r<i,d,r
```

```
        A[j][k-1] + d, //deletion
```

```
        A[j-1][k] + i, //Insertion
```

```
        A[j-1][k-1] + equal(f[k-1],t[j-1])); //Replace or match
```

`equal()` returns m if characters match, r otherwise

Final score is at $A[|t|][|f|]$

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

In action: from erat to arts, Match (+1), Insert/Delete/Replace (-1)

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

r					
t	-3	-3	-1	-1	0
s	-4	-4	2	-2	-1

Global Edit Distance: -1 (Replace, Match, Delete, Match, Insert)

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

Algorithm actually depends on parameter!

<https://eduassistpro.github.io>

```
A[j-1][k-1] + equal(f[k-1],t[j-1])); //Replace or match
```

→ Match score greater than Insert/Del

Add WeChat [edu_assist_pr](#)

e.g. Match (+1), Insert/Delete/Replace (-)

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

Algorithm actually depends on parameter!

<https://eduassistpro.github.io>

→ Match score less than Insert/Delete

Add WeChat edu_assist_pro

(Levenshtein Distance)

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

lengths, e.g. a word and a sentence

Add WeChat edu_assist_pr

Local Edit Distance Algorithm

From string f to string t , given array A of $|f| + 1$ columns and $|t| + 1$ rows, we can solve using the Smith–Waterman algorithm:

```
lf = strlen(f), lt = strlen(t);
A[0][0]=0;
```

```
for (k=1; k<=lf; k++)
```

```
    A[j][k] = max4( //Or min4 if m<i,d,r
```

```
        0,
        A[j][k-1] + d, //deletion
```

```
        A[j-1][k] + i, //Insertion
```

```
        A[j-1][k-1] + equal(f[k-1],t[j-1])); //Replace or match
```

`equal()` returns m if characters match, r otherwise

Final score is greatest value in the entire table (or least value, if $m < i, d, r$)

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

In action: from cart to arts, Match (+1), Insert/Delete/Replace (-1)

r					
t	0	0	0	1	3
s	0	0	0	0	2

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Best match: art with art (+3); ties are possible.

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

For strings f and t , Both algorithms above are $\mathcal{O}(|f||t|)$ in both space and time. (Space can be improved, but time (probably) cannot.)

<https://eduassistpro.github.io> ^{high we want}

$$\mathcal{O}(|f| \sum_{t \in D} |t|)$$

Hence, integer comparisons are roughly t
dictionary. Whether this is feasible depends on t .

Add WeChat edu_assist_pro

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

N-Gram Distance has same goal as Edit Distance: compare two strings to determine “best” match

<https://eduassistpro.github.io>

2-grams of cart: #c, ca, ar, rt, t#

2-grams of arts: #a, ar, rt, ts, s#

Add WeChat [edu_assist_pr](https://eduassistpro.github.io)

N-Gram Distance between n -grams (s) and (t) :

$$|G_n(s)| + |G_n(t)| - 2 \times |G_n(s) \cap G_n(t)|$$

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

$$|G_2(\text{crat})| + |G_2(\text{cart})| - 2 \times |G_2(\text{cart})|$$

$$= 5 + 5 - 2 \times 2 = 6 \text{ (better)}$$

2-Gram Distance between crat and

$$|G_2(\text{crat})| + |G_2(\text{arts})| - 2 \times |G_2(\text{arts})|$$

$$= 5 + 5 - 2 \times 0 = 10$$

Add WeChat [edu_assist_pro](https://eduassistpro.github.io)

Approximate
String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

Occasionally useful as a simpler variant of Edit Distance

<https://eduassistpro.github.io>

Despite its simplicity, takes roughly the same time to compare entire dictionary

Quite useless for very long strings and/or very large dictionaries

Add WeChat [edu_assist_pro](#)

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact
Approximate
Application

Methods

Neighbourhood
Edit Distance
N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

In English (and some other languages), **orthography** (spelling) isn't a

<https://eduassistpro.github.io>

George O'Connell

Also relevant in spelling correction (English)
(so really!)

Add WeChat edu_assist_pro

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

One mechanism: Soundex

a b c d e f g h i j k l m n o p q r s t u v w x y z	→	0 (vowels)
	→	1 (labials)
		2 (misc: fricatives, velars, etc.)

<https://eduassistpro.github.io>

Four step process:

- 1 Except for initial character, translate table
- 2 Remove duplicates (e.g. 4444 → 4)
- 3 Remove 0s
- 4 Truncate to four symbols

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

One mechanism: Soundex

aehiouwy → 0 (vowels)
bpfv → 1 (labials)

tc.)

<https://eduassistpro.github.io>

r → 6 (rhotic)

Four step process:

king	kyngge
k052	k05220
k052	k0520
k52	k52

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

Better phonetic methods make use of the fact that some letters sounds

<https://eduassistpro.github.io>

lpadist uses a text-to-sound algorithm t
the International Phonetic Alphabet (but c

Add WeChat edu_assist_pr

There are also worse variants, like Phonix.

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

Evaluation: consider whether the system is effective at solving the user's problem

<https://eduassistpro.github.io>

To evaluate, we need:

- A number of cases of misspelled word
- The intended (correct) word for each
- An **evaluation metric**

Add WeChat edu_assist_pro

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

We have some cases:

			ght/Wrong?
cracheyt	crotchety	cachet	✓
...	...		

Add WeChat edu_assist_pr

Accuracy: fraction of correct responses (-

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

More realistic situation:

Misspelled Word	Correct Word	Predicted Word	Right/Wrong?
			×
			✓
		carrier	✓
cracheyt	crotchety		

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Precision: fraction of correct responses among attempted responses
($\frac{2}{5}$)

Recall: proportion of words with a correct response (somewhere) ($\frac{2}{3}$)

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Typically, the value of the evaluation metric has little intrinsic meaning

“This system gets 81% accuracy” — useful for users, or not?

<https://eduassistpro.github.io>

the accuracy becomes 74%”

“System A gets 45% precision and 80% recall

System B gets 95% precision and 10% recall

— Which one should we use? (Also: why?)

The answer depends on the problem (and the user)!

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

- What is approximate string search?

■ <https://eduassistpro.github.io>

string? What do we need to generate them?

- How can we evaluate a typical approxi

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Needleman, Saul B. and Wunsch, Christian D. (1970). "A general method applicable to the search for similarities in the amino acid sequence of two proteins". *Journal of Molecular Biology* 48 (3): 443-53.

doi:10.1016/0022-2836(70)90057-4

(Originally in Russian, published in English as:) Levenshtein, Vladimir I.

<https://eduassistpro.github.io>

47:

195197. doi:10.1016/0022-2836(81)9

Kondrak, Grzegorz (2005). "N-Gram String Similarity". *Proceedings of the 12th international conference on Information Retrieval (SPIRE'05)*, pp. 115-126, Buenos Aires, Argentina.

Zobel, Justin and Dart, Philip (1996). "Phonetic String Matching: Lessons from Information Retrieval". In *Proceedings of the 19th annual international ACM SIGIR conference on Research and development in information retrieval (SIGIR'96)*, pp. 166-172, New York, USA.

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

Whitelaw, Casey and Hutchison, Ben and Chung, Grace Y and Ellis, Gerard (2009). "Using the Web for Language Independent

<https://eduassistpro.github.io>

Ahmad, Farooq and Kondrak, Grzegorz (2005). "Learning a Spelling Error Model from Search Query Logs". In Proc

Add WeChat: [edu_assist_pro](#)

Technology Conference and Conference
Natural Language Processing (HLT/EMNLP)
Vancouver, Canada.

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>
g)

- Possibly with “errors” (nucleotide/a

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Typical Genomics problem:

■ <https://eduassistpro.github.io>

- But **much** larger strings: a small genome comparing perhaps 1K character sequences; alphabet is small

Add WeChat edu_assist_pro

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

Recall: we have a “short” (1K character) nucleotide/amino acid

<https://eduassistpro.github.io>

sequence of interest, they might be susceptible to some medical condition

Add WeChat [edu_assist_pr](#)

We're allowed ~10 errors; alphabet s

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

... Forget it.

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

→ Prefers shorter chromosomes (not int

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

... Seems like the right idea.

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

... Can't fit table into memory.

... Requires approximate solutions with he

Add WeChat edu_assist_pr

Approximate String Matching

COMP90049
Knowledge
Technologies

String Search

Exact

Approximate

Application

Methods

Neighbourhood

Edit Distance

N-Gram Distance

Phonetics

Evaluation

References

Genomics

Assignment Project Exam Help

<https://eduassistpro.github.io>

But better methods for using n -grams

Add WeChat edu_assist_pro