

**Approximate
String Search
and Matching**

COMP90049
COMP30018
Knowledge
Technologies

Approximate String Search and Matching

**COMP90049
COMP30018
Knowledge Technologies**

Lea Frermann and Justin Zobel and Karin Verspoor

Semester 2, 2019



THE UNIVERSITY OF

MELBOURNE

String Search
Exact
Approximate
Application
Text
Pre-processing
Methods
Neighbourhood
Edit Distance
References
Phonetic matching (if
there's time)

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Last week... Probability and similarity

Nuts and bolts of knowledge technologies!

- quantifying similarity of complex structures
- estimate the (conditional, joint) probability of observations
- identify high entropy (=informative) information

This week... Approximate String Search and Matching

Another bolt!

- back to similarity
- methods, applications, evaluation

Exact String Search

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

- Given a string, is some substring contained within it?
- Given a string (document), find all occurrences of some substring

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search Exact

Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

- Given a string, is some substring contained within it?
- Given a string (document), find all occurrences of some substring

For example, find Exxon in:

In exes for foxes rex dux mixes a pox of waxed luxes.

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

Grexit or Brexit as quixotic haxxers with buxom rex taxation.

Exact String Search

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search Exact

Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

- Given a string, is some substring contained within it?
- Given a string (document), find all occurrences of some substring

For example, find Exxon in:

In exes for foxes rex dux mixes a pox of waxed luxes.

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

Grexit or Brexit as quixotic haxxers with buxom rex taxation.

Exact String Search

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search Exact

Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

- Given a string, is some substring contained within it?
- Given a string (document), find all occurrences of some substring

For example, find Exxon in:

In exes for foxes rex dux mixes a pox of waxed luxes.

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

Grexit or Brexit as quixotic haxxers with buxom rex taxation.

Not (really) a Knowledge Technology!

Approximate String Search

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Find exon in:

In exes for foxes rex dux mixes a pox of waxed luxes.

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

Grexit or Brexit as quixotic haxxers with buxom rex taxation.

Approximate String Search

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Find `exon` in:

In `exes` for `foxes` `rex` `dux` mixes a `pox` of `waxed` `luxes`.

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

`Grexit` or `Brexit` as `quixotic` `haxxers` with `buxom` `rex` `taxation`.

Not present!

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

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Find `exon` in:

In `exes` for `foxes` `rex` `dux` `mixes` a `pox` of `waxed` `luxes`.

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

`Grexit` or `Brexit` as `quixotic` `haxxers` with `buxom` `rex` `taxation`.

Not present!

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

This is a Knowledge Technology!

Spelling Correction

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)



Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Need the notion of a **dictionary**:

- Here, a list of words

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Need the notion of a **dictionary**:

- Here, a list of ~~words~~ entries that are “correct” with respect to our (expectations of our) language

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Need the notion of a **dictionary**:

- Here, a list of ~~words~~ entries that are “correct” with respect to our (expectations of our) language
- We can break our input into ~~words~~ substrings that we wish to match, and compare each of them against the entries in the dictionary

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Need the notion of a **dictionary**:

- Here, a list of ~~words~~ entries that are “correct” with respect to our (expectations of our) language
- We can break our input into ~~words~~ substrings that we wish to match, and compare each of them against the entries in the dictionary
- A ~~word~~ item in the input which *doesn't* appear in the dictionary is *misspelled*

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Need the notion of a **dictionary**:

- Here, a list of ~~words~~ entries that are “correct” with respect to our (expectations of our) language
- We can break our input into ~~words~~ substrings that we wish to match, and compare each of them against the entries in the dictionary
- A ~~word~~ item in the input which *doesn't* appear in the dictionary is *misspelled*
- A ~~word~~ item in the input which *does* appear in the dictionary might be correctly spelled *or* misspelled (probably slightly beyond the scope of this subject)

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Therefore, the problem here:

Given some item of interest — which does not appear in our dictionary
— which entry from the dictionary was truly intended?

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Therefore, the problem here:

Given some item of interest — which does not appear in our dictionary
— which entry from the dictionary was truly intended?

Depends on the person who wrote the original string!

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Word Blending:

- forming novel words by “blending” two other words
- not simple concatenation (e.g., football is not a blend word)

breakfast	+	lunch	→	brunch
fork	+	spoon	→	spork
Britain	+	exit	→	Brexit

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Word Blending:

- forming novel words by “blending” two other words
- not simple concatenation (e.g., football is not a blend word)

breakfast	+	lunch	→	brunch
fork	+	spoon	→	spork
Britain	+	exit	→	Brexit

- Language changes constantly!
- New terms often coined in colloquial language (e.g., Twitter)
- Can we build **knowledge technologies** that detect novel blends?
...or help us understand what their **components** are?
- **Assignment 1**

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Name matching

the name *Gorbachev* is spelled (at least) 20 different ways in a corpus of newswire text!

Gorbachev, Gorbacahev, Gorbahev, Gorbatchev, Gorbechev,
Gorbachov, Gorachev, Gorbacheva, Gorbechyev, Gorbacev,
Gorbachyov, Gorabchev, Grobachev, ...

street and place name conventions

boulevard|blvd|bd|bde|blv|bl|blvde|blvrd|boulavard|boul|bvd
apartment|apt|ap|aprt|apmnt
village|vil|vge|vill|villag|villg|vlg|vlge|vllg

Other Problems of Interest

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

- Computational Genomics
- Name matching
- Query repair
- Phonetic matching
- Data cleaning (e.g., deduplication)
- ...

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Before we start... fix some irregularities which may distract the matching algorithm.

Aside: Text Normalization

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Original: I had an AMAZING trip to Italy, Coffee
is only 2 bucks, sometimes three! it's
incredible!!! Right?

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Original: I had an AMAZING trip to Italy, Coffee
is only 2 bucks, sometimes three! it's
incredible!!! Right?

Normalized: i had an amazing trip to italy, coffee
is only 2 bucks, sometimes three! it's
incredible!!! right?

1. Consistent Casing

map all characters to upper (or lower) case

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Original: I had an AMAZING trip to Italy, Coffee
is only 2 bucks, sometimes three! it's
incredible!!! Right?

Normalized: i had an amazing trip to italy, coffee
is only two bucks, sometimes three! it's
incredible!!! right?

2. Unify or remove numbers

map numbers to a consistent representation or remove (replace) them
entirely

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Original: I had an AMAZING trip to Italy, Coffee
is only 2 bucks, sometimes three! it's
incredible!!! Right?

Normalized: i had an amazing trip to italy, coffee
is only two bucks, sometimes three! it's
incredible!!! right?

3. Remove unnecessary spacing

e.g., duplicate white spaces

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Original: I had an AMAZING trip to Italy, Coffee
is only 2 bucks, sometimes three! it's
incredible!!! Right?

Normalized: i had an amazing trip to italy , coffee is
only two bucks , sometimes three ! it 's
incredible ! ! ! Right ?

4. Word Tokenization

split text into individual tokens (aka words)

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Original: I had an AMAZING trip to Italy, Coffee
is only 2 bucks, sometimes three! it's
incredible!!! Right?

Normalized: i had an amazing trip to italy , coffee is
only two bucks , sometimes three ! it 's
incredible ! ! ! Right ?

4. Word Tokenization

split text into individual tokens (aka words)

What's a “best” match?

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Find approximate match(es) for `exon` in:

In exes for foxes rex dux mixes a pox of waxed luxes.

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

Grexit or Brexit as quixotic haxxers with buxom rex taxation.

What's a “best” match?

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Find approximate match(es) for `exon` in:

In exes for foxes rex dux mixes a pox of waxed luxes.

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

Grexit or Brexit as quixotic haxxers with buxom rex taxation.

`exon` → `Exxon` **Insert x**

`exon` → `exo` **Delete n**

`exon` → `axon` **Replace e with a (Sometimes Substitute)**

`exon` → `oxen` **Transpose e and o (not covered here)**

What's a "best" match?

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Find approximate match(es) for `exon` in:

In exes for foxes rex dux mixes a pox of waxed luxes.

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

Grexit or Brexit as quixotic haxxers with buxom rex taxation.

`exon` → `Exxon` **Insert** x

`exon` → `exo` **Delete** n

`exon` → `axon` **Replace** e with a (Sometimes **Substitute**)

`exon` → `oxen` **Transpose** e and o (not covered here)

What's a "best" match?

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Find approximate match(es) for `exon` in:

In exes for foxes rex dux mixes a pox of waxed luxes.

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

Grexit or Brexit as quixotic haxxers with buxom rex taxation.

`exon` → `Exxon` **Insert** x

`exon` → `exo` **Delete** n

`exon` → `axon` **Replace** e with a (Sometimes **Substitute**)

`exon` → `oxen` **Transpose** e and o (not covered here)

What's a "best" match?

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Find approximate match(es) for `exon` in:

In exes for foxes rex dux mixes a pox of waxed luxes.

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

Grexit or Brexit as quixotic haxxers with buxom rex taxation.

`exon` → Exxon **Insert x**

`exon` → exO **Delete n**

`exon` → axon **Replace e with a (Sometimes Substitute)**

`exon` → oxen **Transpose e and o (not covered here)**

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

For a given string w of interest:

- Generate all variants of w that utilise at most k changes (Insertions/Deletions/Replacements) — **neighbours**

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

For a given string w of interest:

- Generate all variants of w that utilise at most k changes (Insertions/Deletions/Replacements) — **neighbours**
- Check whether generated variants exist in dictionary

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

For a given string w of interest:

- Generate all variants of w that utilise at most k changes (Insertions/Deletions/Replacements) — **neighbours**
- Check whether generated variants exist in dictionary
- **All** results found in dictionary are returned

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

For a given string w of interest:

- Generate all variants of w that utilise at most k changes (Insertions/Deletions/Replacements) — **neighbours**
- Check whether generated variants exist in dictionary
- **All** results found in dictionary are returned

Unix command-line utility `agrep` is an efficient tool for finding these.

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

For a given string w of interest:

- Generate all variants of w that utilise at most k changes (Insertions/Deletions/Replacements) — **neighbours**
- Check whether generated variants exist in dictionary
- **All** results found in dictionary are returned

Unix command-line utility `agrep` is an efficient tool for finding these.

For example:

... proceed if you can see no **ther** option ...

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

For a given string w of interest:

- Generate all variants of w that utilise at most k changes (Insertions/Deletions/Replacements) — **neighbours**
- Check whether generated variants exist in dictionary
- **All** results found in dictionary are returned

Unix command-line utility `agrep` is an efficient tool for finding these.

For example:

... proceed if you can see no **ther** option ...

the their there tier **other** mther tnher thpr ...

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Neighborhood search

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Neighborhood search

Consider: alphabet size is Σ , length of string is $|w|$:

For 1 edit, roughly $\mathcal{O}(\Sigma \cdot |w|)$ neighbours

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Neighborhood search

Consider: alphabet size is Σ , length of string is $|w|$:

For 2 edits, roughly $\mathcal{O}(\Sigma^2 \cdot |w|^2)$ neighbours

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Neighborhood search

Consider: alphabet size is Σ , length of string is $|w|$:

For k edits, roughly $\mathcal{O}(\Sigma^k \cdot |w|^k)$ neighbours

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Neighborhood search

Consider: alphabet size is Σ , length of string is $|w|$:

For k edits, roughly $\mathcal{O}(\Sigma^k \cdot |w|^k)$ neighbours

But,

- Σ is a small constant
- string of interest is usually short
- k is usually small

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Neighborhood search

Consider: alphabet size is Σ , length of string is $|w|$:

For k edits, roughly $\mathcal{O}(\Sigma^k \cdot |w|^k)$ neighbours

But,

- Σ is a small constant
- string of interest is usually short
- k is usually small

Dictionary Read

Assuming D entries, binary search yields $\mathcal{O}(|w|^k \log D)$ string comparisons

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Neighborhood search

Consider: alphabet size is Σ , length of string is $|w|$:

For k edits, roughly $\mathcal{O}(\Sigma^k \cdot |w|^k)$ neighbours

But,

- Σ is a small constant
- string of interest is usually short
- k is usually small

agrep example

Dictionary Read

Assuming D entries, binary search yields $\mathcal{O}(|w|^k \log D)$ string comparisons

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

An alternative method:

Scan through each dictionary entry looking for the “best” match

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

An alternative method:

Scan through each dictionary entry looking for the “best” match

Intuition:

- Transform the string of interest into each dictionary entry
- **Operations:** Insert, Delete, Replace, and Match

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

An alternative method:

Scan through each dictionary entry looking for the “best” match

Intuition:

- Transform the string of interest into each dictionary entry
- **Operations:** Insert, Delete, Replace, and Match
- Each operation is associated with a score;

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

An alternative method:

Scan through each dictionary entry looking for the “best” match

Intuition:

- Transform the string of interest into each dictionary entry
- **Operations:** Insert, Delete, Replace, and Match
- Each operation is associated with a score;
- Best match is the dictionary entry with best aggregate **score**

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

An alternative method:

Scan through each dictionary entry looking for the “best” match

Intuition:

- Transform the string of interest into each dictionary entry
- **Operations:** Insert, Delete, Replace, and Match
- Each operation is associated with a score;
- Best match is the dictionary entry with best aggregate **score**

Global Edit Distance – Example

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Item of interest	crat
Dictionary	cart, arts
Scores	Match +1, Insert -1, Delete -1, Replace -1

1 Transform crat → cart

Global Edit Distance – Example

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Item of interest	crat
Dictionary	cart, arts
Scores	Match +1, Insert -1, Delete -1, Replace -1

1 Transform crat → cart

Match c Delete r Match a Insert r Match t

Global Edit Distance – Example

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Item of interest `crat`

Dictionary `cart, arts`

Scores Match +1, Insert -1, Delete -1, Replace -1

1 Transform `crat` → `cart`

Match c	Delete r	Match a	Insert r	Match t	
+1	-1	+1	-1	+1	= +1

Global Edit Distance – Example

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Item of interest `crat`

Dictionary `cart, arts`

Scores Match +1, Insert -1, Delete -1, Replace -1

1 Transform `crat` → `cart`

Match c	Delete r	Match a	Insert r	Match t	
+1	-1	+1	-1	+1	= +1

2 Transform `crat` → `arts`

Global Edit Distance – Example

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Item of interest	crat
Dictionary	cart, arts
Scores	Match +1, Insert -1, Delete -1, Replace -1

1 Transform crat → cart

Match c	Delete r	Match a	Insert r	Match t	
+1	-1	+1	-1	+1	= +1

2 Transform crat → arts

Replace c, a	Match r	Delete a	Match t	Insert s
--------------	---------	----------	---------	----------

Global Edit Distance – Example

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Item of interest crat

Dictionary cart, arts

Scores Match +1, Insert -1, Delete -1, Replace -1

1 Transform crat → cart

Match c	Delete r	Match a	Insert r	Match t	
+1	-1	+1	-1	+1	= +1

2 Transform crat → arts

Replace c, a	Match r	Delete a	Match t	Insert s	
-1	+1	-1	+1	-1	= -1

Global Edit Distance – Example

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Item of interest crat

Dictionary cart, arts

Scores Match +1, Insert -1, Delete -1, Replace -1

1 Transform crat → cart

Match c	Delete r	Match a	Insert r	Match t	
+1	-1	+1	-1	+1	= +1

2 Transform crat → arts

Replace c, a	Match r	Delete a	Match t	Insert s	
-1	+1	-1	+1	-1	= -1

cart is the better match

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Confusingly, Global Edit Distance isn't necessarily a “distance”

Why?

Confusingly, Global Edit Distance isn't necessarily a “distance”

Why?

But we can make it one

- Global edit distance with parameters:

Match (0), Insert (+1), Delete (+1), Replace (+1)

- counts the number of edits required to transform one string into the other
- **Levenshtein Distance**

Approximate
String Search
and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search
Exact

Approximate
Application

Text
Pre-processing

Methods
Neighbourhood
Edit Distance

References
Phonetic matching (if
there's time)

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Confusingly, Global Edit Distance isn't necessarily a “distance”

Why?

But we can make it one

- Global edit distance with parameters:

Match (0), Insert (+1), Delete (+1), Replace (+1)

- counts the number of edits required to transform one string into the other

- **Levenshtein Distance**

Often,

- $\text{cost}(\text{Insert}) = \text{cost}(\text{Delete})$ (direction doesn't matter)

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Confusingly, Global Edit Distance isn't necessarily a “distance”

Why?

But we can make it one

- Global edit distance with parameters:

Match (0), Insert (+1), Delete (+1), Replace (+1)

- counts the number of edits required to transform one string into the other

- **Levenshtein Distance**

Often,

- $\text{cost}(\text{Insert}) = \text{cost}(\text{Delete})$ (direction doesn't matter)
- score of Replace depends on the character being replaced

Is faxing more likely to be facing or faking?

Global Edit Distance – Example

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Item of interest	crat
Dictionary	cart, arts
Scores	Match 0, Insert +1, Delete +1, Replace +1

1 Transform crat → cart

Match c	Delete r	Match a	Insert r	Match t	
0	+1	0	+1	0	= +2

2 Transform crat → arts

Replace c, a	Match r	Delete a	Match t	Insert s	
+1	0	+1	0	+1	= +3

3 changes vs 2 changes: cart is the better match

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Hypothetically, any parameter is possible!

But some choices make no sense, e.g.:

Match (+4), Insert (-2), Delete (+8), Replace (0)

Consider aba: which corresponds to best match?

- foo
- aba
- cbc

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Hypothetically, any parameter is possible!

But some choices make no sense, e.g.:

Match (+4), Insert (-2), Delete (+8), Replace (0)

Consider aba: which corresponds to best match?

- foo: Insert, Delete, Insert, Delete, Insert, Delete
- aba: Match, Match, Match
- cbc: Replace, Match, Replace

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Hypothetically, any parameter is possible!

But some choices make no sense, e.g.:

Match (+4), Insert (-2), Delete (+8), Replace (0)

Consider aba: which corresponds to best match?

- foo: Insert, Delete, Insert, Delete, Insert, Delete = +18
- aba: Match, Match, Match = +12
- cbc: Replace, Match, Replace = +4

Computer can't find best sequence of operations by inspection

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Global Edit Distance Algorithm

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

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 (j=1; j<=lt; j++) A[j][0] = j * i;
for (k=1; k<=lf; k++) A[0][k] = k * d;

for (j=1; j<=lt; j++)
    for (k=1; k<=lf; k++)
        A[j][k] = max3( //Or min3 if m<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 $A[|t|][|f|]$

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

	ϵ	c	r	a	t
ϵ					
a					
r					
t					
s					

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

Initialise table:

	ϵ	c	r	a	t
ϵ	0	-1	-2	-3	-4
a	-1				
r	-2				
t	-3				
s	-4				

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

For `c`–`a` correspondence, consider three neighbours:

	ϵ	c	r	a	t
ϵ	0	-1	-2	-3	-4
a	-1	?			
r	-2				
t	-3				
s	-4				

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

For `c`—`a` correspondence, Delete `c`:

	ϵ	c	r	a	t
ϵ	0	-1	-2	-3	-4
a	-1	-2			
r	-2				
t	-3				
s	-4				

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

For c-a correspondence, Insert a:

	ϵ	c	r	a	t
ϵ	0	-1	-2	-3	-4
a	-1	-2			
r	-2				
t	-3				
s	-4				

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

For c—a correspondence, Replace c with a:

	ϵ	c	r	a	t
ϵ	0	-1	-2	-3	-4
a	-1	-1			
r	-2				
t	-3				
s	-4				

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

And so on:

	ϵ	c	r	a	t
ϵ	0	-1	-2	-3	-4
a	-1	-1	-2		
r	-2				
t	-3				
s	-4				

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

And so on:

	ε	c	r	a	t
ε	0	-1	-2	-3	-4
a	-1	-1	-2	-1	
r	-2				
t	-3				
s	-4				

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods
Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

And so on:

	ε	c	r	a	t
ε	0	-1	-2	-3	-4
a	-1	-1	-2	-1	-2
r	-2	-2	0	-1	-2
t	-3	-3	-1	-1	0
s	-4	-4	-2	-2	-1

Global Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

And so on:

	ε	c	r	a	t
ε	0	-1	-2	-3	-4
a	-1	-1	-2	-1	-2
r	-2	-2	0	-1	-2
t	-3	-3	-1	-1	0
s	-4	-4	-2	-2	-1

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

More parameter concerns

Algorithm actually depends on parameter!

Approximate
String Search
and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

More parameter concerns

Algorithm actually depends on parameter!

```
A[j][k] = max3(  
    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
```

Approximate
String Search
and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search
Exact

Approximate
Application

Text
Pre-processing

Methods
Neighbourhood
Edit Distance

References
Phonetic matching (if
there's time)

More parameter concerns

Algorithm actually depends on parameter!

```
A[j][k] = max3(  
    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
```

→ Match score greater than Insert/Delete/Replace

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

Approximate
String Search
and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search
Exact

Approximate
Application

Text
Pre-processing

Methods
Neighbourhood
Edit Distance

References
Phonetic matching (if
there's time)

More parameter concerns

Algorithm actually depends on parameter!

```
A[j][k] = min3(  
    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
```

Approximate
String Search
and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search
Exact

Approximate
Application

Text
Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

More parameter concerns

Algorithm actually depends on parameter!

```
A[j][k] = min3(  
    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
```

→ Match score less than Insert/Delete/Replace

e.g. Match (0), Insert/Delete/Replace (+1)

(Levenshtein Distance)

Approximate
String Search
and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search
Exact

Approximate
Application

Text
Pre-processing

Methods
Neighbourhood
Edit Distance

References
Phonetic matching (if
there's time)

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

like Global Edit Distance, but we are searching for the best
substring match

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

**like Global Edit Distance, but we are searching for the best
substring match**

Particularly suitable when comparing two strings of very different
lengths, e.g.

- a suffix of a word
- a word in a sentence
- a sentence in an entire document

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 (j=1; j<=lt; j++) A[j][0] = 0;
for (k=1; k<=lf; k++) A[0][k] = 0;

for (j=1; j<=lt; j++)
    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 Search
and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

(For Local Edit Distance, Match must have different $+/-$ sign to Insert/Delete/Replace)

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

	ϵ	c	a	r	t
ϵ					
a					
r					
t					
s					

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

Initialise table:

	ϵ	c	a	r	t
ϵ	0	0	0	0	0
a	0				
r	0				
t	0				
s	0				

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

For c-a correspondence, consider three neighbours:

	ϵ	c	a	r	t
ϵ	0	0	0	0	0
a	0	?			
r	0				
t	0				
s	0				

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

For c-a correspondence, Delete c:

	ϵ	c	a	r	t
ϵ	0	0	0	0	0
a	0	-1			
r	0				
t	0				
s	0				

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

For c-a correspondence, Insert a:

	ϵ	c	a	r	t
ϵ	0	0	0	0	0
a	0	-1			
r	0				
t	0				
s	0				

Local Edit Distance in Action

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

→ delete source word

↓ insert target word

For c-a correspondence, Replace c with a:

	ϵ	c	a	r	t
ϵ	0	0	0	0	0
a	0	-1			
r	0				
t	0				
s	0				

Approximate
String Search
and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

For c-a correspondence, 0 is better:

	ϵ	c	a	r	t
ϵ	0	0	0	0	0
a	0	0			
r	0				
t	0				
s	0				

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

For a—a correspondence (Match), 1 is better:

	ϵ	c	a	r	t
ϵ	0	0	0	0	0
a	0	0	1		
r	0				
t	0				
s	0				

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

For a-r correspondence, back to 0:

	ϵ	c	a	r	t
ϵ	0	0	0	0	0
a	0	0	1	0	
r	0				
t	0				
s	0				

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

And so on:

	ϵ	c	a	r	t
ϵ	0	0	0	0	0
a	0	0	1	0	0
r	0	0	0	2	1
t	0	0	0	1	3
s	0	0	0	0	2

Local Edit Distance in Action

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

→ delete source word

↓ insert target word

And so on:

	ϵ	c	a	r	t
ϵ	0	0	0	0	0
a	0	0	1	0	0
r	0	0	0	2	1
t	0	0	0	1	3
s	0	0	0	0	2

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

Edit Distance Efficiency

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

- single source string f
- multiple target strings t in dictionary D
- a single approximate match (both global and local)

$$\mathcal{O}(|f||t|)$$

Edit Distance Efficiency

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

- single source string f
- multiple target strings t in dictionary D
- a single approximate match (both global and local)

$$\mathcal{O}(|f||t|)$$

- approximate matching of f over the whole dictionary

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

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

- single source string f
- multiple target strings t in dictionary D
- a single approximate match (both global and local)

$$\mathcal{O}(|f||t|)$$

- approximate matching of f over the whole dictionary

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

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

- single source string f
- multiple target strings t in dictionary D
- a single approximate match (both global and local)

$$\mathcal{O}(|f||t|)$$

- approximate matching of f over the whole dictionary

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

- **feasibility depends (linearly) on size of the dictionary**

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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. (1966). "Binary codes capable of correcting deletions, insertions, and reversals". *Soviet Physics Doklady* 10 (8): 707–710.

Christin, P. (2006). "A Comparison of Personal Name Matching: Techniques and Practical Issues". *Proceedings of the Sixth IEEE International Conference on Data Mining - Workshops*

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

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

Peng, N. and Yu, M. and Drezde, M. (2015). "An Empirical Study of Chinese Name Matching and Applications". In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pp. 377–383, Beijing, China.

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

Whitelaw, Casey and Hutchison, Ben and Chung, Grace Y and Ellis, Gerard (2009). "Using the Web for Language Independent Spellchecking and Autocorrection". In Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing (EMNLP 2009), pp. 890-899, Singapore, Singapore.

Ahmad, Farooq and Kondrak, Grzegorz (2005). "Learning a Spelling Error Model from Search Query Logs". In Proceedings of the Human Technology Conference and Conference on Empirical Methods in Natural Language Processing (HLT/EMNLP 2005), pp. 955-962, Vancouver, Canada.

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

One (ineffectual) mechanism: Soundex

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

One mechanism: Soundex

Translation table:	aehiouwy	→	0 (vowels)
	bpfv	→	1 (labials)
	cgjksxz	→	2 (misc: fricatives, velars, etc.)
	dt	→	3 (dentals)
	l	→	4 (lateral)
	mn	→	5 (nasals)
	r	→	6 (rhotic)

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

One mechanism: Soundex

	aehiouwy	→	0 (vowels)
	bpfv	→	1 (labials)
	cgjkqszx	→	2 (misc: fricatives, velars, etc.)
Translation table:	dt	→	3 (dentals)
	l	→	4 (lateral)
	mn	→	5 (nasals)
	r	→	6 (rhotic)

Four step process:

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

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

One mechanism: Soundex

Translation table:	aehiouwy	→	0 (vowels)
	bpfv	→	1 (labials)
	cgjksxz	→	2 (misc: fricatives, velars, etc.)
	dt	→	3 (dentals)
	l	→	4 (lateral)
	mn	→	5 (nasals)
	r	→	6 (rhotic)

Four step process:

king	kyngge
k052	k05220
k052	k0520
k52	k52

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

One mechanism: Soundex

Translation table:	aehiouwy	→	0 (vowels)
	bpfv	→	1 (labials)
	cgjksxz	→	2 (misc: fricatives, velars, etc.)
	dt	→	3 (dentals)
	l	→	4 (lateral)
	mn	→	5 (nasals)
	r	→	6 (rhotic)

Four step process:

knight	night
k50203	n0203
k50203	n0203
k523	n23

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact

Approximate

Application

Text

Pre-processing

Methods

Neighbourhood

Edit Distance

References

Phonetic matching (if
there's time)

One mechanism: Soundex

Translation table:	aehiouwy	→	0 (vowels)
	bpfv	→	1 (labials)
	cgjksxz	→	2 (misc: fricatives, velars, etc.)
	dt	→	3 (dentals)
	l	→	4 (lateral)
	mn	→	5 (nasals)
	r	→	6 (rhotic)

Four step process:

loan	loew	lough	lewicks
1005	1000	10020	1000222
105	10	1020	102
15	1	12	12

Approximate String Search and Matching

COMP90049
COMP30018
Knowledge
Technologies

String Search

Exact
Approximate
Application

Text

Pre-processing

Methods

Neighbourhood
Edit Distance

References

Phonetic matching (if
there's time)

Better phonetic methods make use of the fact that some letters sounds alike in certain contexts, and different in other contexts

Editex uses the Edit Distance to compare strings based on a similar translation table to Soundex

Ipadist uses a text-to-sound algorithm to represent tokens according to the International Phonetic Alphabet (but context matters a lot)

There are also worse variants, like Phonix.