Stemming algorithms

Modern Theoretical Computer Science 2014 / 2015

Author:Ing. Petr LoukotaAdvisor:Doc. Ing. Jaroslav Zendulka, CSc.

Contents

- Search Engines and Web Page Classification
- Stemming, difference between Stemming and Lemmatization
- Classification of Stemming Algorithms
- Truncating Methods, Statistical Methods, Inflectional and Derivational Methods, Corpus Based Methods and Context Sensitive Methods
- Porter Stemming Algorithm
- Errors in Stemming

How Search Engines work



Fakulta informačních technologií VUT v Brně

Internet Mapy Obrázky Zprávy Vídea Více * Vyhledávací nástroje



Přibližný počet výsledků: 98 900 (0,44 s)

Fakulta informačních technologií VUT v Brně www.fit.vutbr.cz/.cs *

Fakulta informačnich technologii VUT v Brně. ... Vysokého učení technického v Brně. Božetěchova 1/2 612 66 Brno, Czech Republic. Tel.: +420 54114-1144 ...

Přijímací řízení

Přijímací řízení. Důležité termíny a předpisy (úřední deska FIT ...

O fakultě

Informace o fakultě. Fakulta informačních technologií (FIT ...

Další výsledky z webu vutbr.cz »

Studium a výuka

Studium a výuka. Studijní programy na FIT · ECTS a

Studium na FIT

Studium na FIT. [study]. Podrobněji: Bakalářský studijní program ...

How Search Engines work



Fakulta informačních technologií VUT v Brně

Internet Mapy Obrázky Zprávy Vídea Více * Vyhledávací nástroje

Přibližný počet výsledků: 98 900 (0,44 s)

Fakulta informačních technologií VUT v Brně www.fit.vutbr.cz/.cs *

Fakulta informačnich technologii VUT v Brně. ... Vysokého učení technického v Brně. Božetěchova 1/2 612 66 Brno, Czech Republic. Tel.: +420 54114-1144 ...

Přijímaci řízení

Přijímací řízení. Důležité termíny a předpisy (úřední deska FIT ...

O fakultě

Informace o fakultě. Fakulta informačních technologií (FIT ...

Další výsledky z webu vutbr.cz »

Studium a výuka

Studium a výuka. Studijní programy na FIT · ECTS a

Studium na FIT

Studium na FIT. [study]. Podrobněji: Bakalářský studijní program ...

How Search Engines work



Fakulta informačních technologií VUT v Brně

Internet Mapy Obrázky Zprávy Videa Více * Vyhledávací nástroje



Přibližný počet výsledků: 98 900 (0,44 s)

Fakulta informačních technologií VUT v Brně www.fit.vutbr.cz/.cs *

Fakulta informačnich technologii VUT v Brně. ... Vysokého učení technického v Brně. Božetěchova 1/2 612 66 Brno, Czech Republic. Tel.: +420 54114-1144 ...

Přijímací řízení

Přijímací řízení. Důležité termíny a předpisy (úřední deska FIT ...

O fakultě

Informace o fakultě. Fakulta informačních technologií (FIT ...

Další výsledky z webu vutbr.cz »

Studium a výuka

Studium a výuka. Studijní programy na FIT · ECTS a

Studium na FIT

Studium na FIT. [study]. Podrobněji: Bakalářský studijní program ...



Web Page Classification

- Classification is a process of assigning some of pre-defined class labels to a web document.
- Classification can take both visual and textual information into consideration.
- Two standard pre-processing procedures are applied on the text obtained from a web page stop words removal and stemming.

Stemming

- Stemming is a pre-processing step in Text Mining applications as well as a very common requirement of Natural Language procession functions.
- It is very important in most of the Information Retrieval systems.
- The main purpose of stemming is to reduce different grammatical forms of a word like its noon, adjective, verb, adverb etc. to its root form.

Stemming

- Stemming is a pre-processing step in Text Mining applications as well as a very common requirement of Natural Language procession functions.
- It is very important in most of the Information Retrieval systems.
- The main purpose of stemming is to reduce different grammatical forms of a word like its noon, adjective, verb, adverb etc. to its root form.



Stemming and Lemmatizing

Stemming:

 In stemming the 'stem' is obtaining after applying a set of rules but without bothering about part of speech or the context of the word occurrence.

Lemmatization:

Lemmatization deals with obtaining the 'lemma' of a word which involves reducing the word forms to its root form after understanding the part of speech and the context of the word in the given sentence.

Stemming and Lemmatizing

Stemming:

 In stemming the 'stem' is obtaining after applying a set of rules but without bothering about part of speech or the context of the word occurrence.



Lemmatization:

Lemmatization deals with obtaining the 'lemma' of a word which involves reducing the word forms to its root form after understanding the part of speech and the context of the word in the given sentence.



Classification of Stemming Algorithms

- Since the meaning is same but the word form is different it is necessary to identify each word form with its base form.
- To do this a variety of Stemming Algorithms have been developed.



Truncating Methods

• These methods are related to removing the suffixes or prefixes (commonly known as affixes) of a word.

Truncating Methods

• These methods are related to removing the suffixes or prefixes (commonly known as affixes) of a word.

Porter Stemming Algorithm:

- It has five steps, and within each step, rules are applied until one of them passes the condition.
- The rule look like the following: <condition> <suffix> -> <new suffix>

Example: (m > 0) EED -> EE

This rule means "if the word has at least one vowel and consonant plus EED ending, change the ending to EE".

So **"agreed**" becomes **"agree**" while **"feed**" remains unchanged.

Statistical Methods

• These are the Stemmers who are based on statistical analysis and techniques. Most of the methods remove the affixes but after implementing some statistical procedure.

Statistical Methods

N-Gram Stemmer:

- String-similarity approach is used to convert word inflation to its stem.
- An n-gram is a set of consecutive characters extracted from a word.
- Similar word will have a high proportion of n-grams in common.

Example: the words "statistics" and "statistical" for n equals to 2 The diagrams: st, ta, at, ti, is, st, ti, ic, cs

st, ta, at, ti, is, st, ti, ic, ca, al

Advantage: It is language independent.

Disadvantage: It requires a significant amount of memory and storage

Statistical Methods

N-Gram Stemmer:

- String-similarity approach is used to convert word inflation to its stem.
- An n-gram is a set of consecutive characters extracted from a word.
- Similar word will have a high proportion of n-grams in common.

Example: the words "statistics" and "statistical" for n equals to 2 The diagrams: st, ta, at, ti, is, st, ti, ic, cs

st, ta, at, ti, is, st, ti, ic, ca, al

Advantage: It is language independent.

Disadvantage: It requires a significant amount of memory and storage

Inflectional and Derivational Methods

- It involves both the inflectional as well as the derivational morphology analysis.
- In case of inflectional the word variants are related to the language specific syntactic variations like plural, gender, case etc. whereas in derivational the world variants are related to the part of speech of a sentence where the word occurs.

Inflectional and Derivational Methods

- It involves both the inflectional as well as the derivational morphology analysis.
- In case of inflectional the word variants are related to the language specific syntactic variations like plural, gender, case etc. whereas in derivational the world variants are related to the part of speech of a sentence where the word occurs.
- Derivational affixes do frequently change the category of the stem they attach to:
 write (verb) -> writer (noun), break (verb) -> breakable (adjective)
 - inflectional affixes never change the category of the stems they attach to:
 - cow (noun) -> cows (noun), hate (verb) -> hated (verb)

Corpus Based Methods

- This approach tries to overcome some of the drawback of Porter Stemmer.
 - For example, the words "policy" and "police" are conflated though they have a different meaning but the words "index" and "indices" are not conflated though they have the same root.
 - Porter Stemmer also generates stems which are not real words like "iteration" becomes "iter" and "general" becomes "gener".
- Corpus based stemming refers to automatic modification of words that have resulted in a common stem to suit the characteristics of a given text corpus using statistical methods.
- Using this concept some of the over-stemming and under-stemming drawbacks are resolved, e.g. "policy" and "police" will no longer be conflated.

Context Sensitive Methods

- This method was proposed by Funchun Peng.
- Unlike the usual method where stemming is done before indexing a document, over here for a Web Search, context sensitive analysis is done using statistical modeling on the query side.
- First, based on statistical language modeling, context sensitive analysis on the query side is performed. This predicts which of its morphological variants is useful to expand a query term with before submitting the query to the search engine.
- Second, context sensitive document matching for those expanded variants is performed.

- One of the most popular Stemming Algorithms proposed in 1980 by Martin Porter
- Many modifications and enhancements have been done and suggested on the basic algorithm.
- It is based on the idea that the suffixes in the English language (approximately 1200) are mostly made up of a combination of smaller and simpler suffixes.
- It has five steps, and within each step, rules are applied until one of them passes the conditions.
- If a rule is accepted, the suffix is removed accordingly.

Consonant – a letter other than A, E, I, O, U, and Y preceded by a consonant
 Example: TOY – consonants are T and Y, SYZYGY – consonants are S, Z, G
 Vowel – a letter that is not a consonant

Consonant will be denoted by c, vowel will be denoted by v. A list ccc of length greater than 0 will be denoted by C. A list vvv of length greater than 0 will be denoted by V.

Any word has one of the four forms: **CVCV** ... **C**, **CVCV** ... **V**, **VCVC** ... **C**, **VCVC** ... **V** -> [**C**]**VCVC** ... [**V**] -> [**C**](**VC**){m}[**V**]

The rules for removing a suffix in the form: (condition) S1 -> S2Example: (m > 1) EMENT ->REPLACEMENT -> REPLAC

The condition may also contain expression with and, or and not and the following:

- ***S** the stem ends with S (similarly for other letters)
- ***v*** the stem contains a vowel
- *d the stem ends with a double consonant (e.g. –TT, -SS)
- *o the stem ends cvc, where the second c is not W, X or Y (e.g. –WIL, -HOP)

Only one rule is obeyed – the one **with the longest matching** S1 for the given word.

a) SSES -> SS	caresses -> caress
IES -> I	pon <mark>ies</mark> -> poni
SS -> SS	caress -> caress
S ->	cats -> cat

b) (m > 0) EED -> EE	agreed-> agree
(*v*) ED ->	plaster <mark>ed</mark> -> plaster
(*v*) ING ->	motoring -> motor

If the second or third of the rules in b) is successful, the following is done: $AT \rightarrow ATE$ $conflat(ed) \rightarrow conflate$ $BL \rightarrow BLE$ $troubl(ed) \rightarrow trouble$ $IZ \rightarrow IZE$ $siz(ed) \rightarrow size$ $(*d and not (*L or *S or *Z)) \rightarrow single letter<math>hopp(ing) \rightarrow hop$ $(m = 1 and *o) \rightarrow E$ $fil(ing) \rightarrow file$

c) (*v*) Y -> I

happy -> happi

• Step 1 deals with plurals and past participles. The subsequent steps are much more straightforward.

(m > 0) ATIONAL -> ATE (m > 0) TIONAL -> TION (m > 0) ENCI -> ENCE (m > 0) ANCI -> ANCE (m > 0) IZER -> IZE (m > 0) ABLI -> ABLE (m > 0) ALLI -> AL (m > 0) ENTLI -> ENT (m > 0) ELI -> E(m > 0) OUSLI -> OUS (m > 0) IZATION -> IZE (m > 0) ATION -> ATE (m > 0) ATOR -> ATE (m > 0) ALISM -> AL

relational -> relate conditional -> condition $valenci \rightarrow valence$ hesitanci -> hesitance digitizer -> digitize conformabli -> conformable radicalli -> radical differentli -> different vileli -> vile analogousli -> analogous vietnamization -> vietnamize predication -> predicate operator -> operate feudalism -> feudal

(m > 0) IVENESS -> IVE
(m > 0) FULNESS -> FUL
(m > 0) OUSNESS -> OUS
(m > 0) ALITI -> AL
(m > 0) IVITI -> IVE
(m > 0) BILITI -> BLE

decisiveness -> decisive hopefulness -> hopeful callousness -> callous formaliti -> formalal sensitiviti -> sensitive sensibiliti -> sensible

(m > 0) ICATE -> IC (m > 0) ATIVE -> (m > 0) ALIZE -> AL (m > 0) ICITI -> IC (m > 0) ICAL -> IC (m > 0) FUL -> (m > 0) NESS -> triplicate -> triplic formative -> form formalize -> formal electriciti -> electric electrical -> electric hopeful -> hope goodness -> good

(m > 1) AL ->(m > 1) ANCE -> (m > 1) ENCE -> (m > 1) ER ->(m > 1) IC -> (m > 1) ABLE -> (m > 1) ANT ->(m > 1) EMENT -> (m > 1) MENT -> (m > 1) ENT -> (m > 1 and (*S or *T)) ION -> (m > 1) OU -> (m > 1) ISM -> (m > 1) ATE ->

revival -> reviv allowance -> allow inference -> infer airliner -> airlin gyroscopic -> gyroscop adjustable -> adjust irritant -> irrit replacement -> replac adjustment -> adjust dependent -> depend adoption -> adopt homologou -> homolog communism -> commun activate -> activ

(m > 1) ITI -> (m > 1) OUS -> (m > 1) IVE -> (m > 1) IZE -> angulariti -> angular homologous -> homolog effective -> effect bowdlerize -> bowdler

• The suffixes are now removed. All that remains is a little tidying up.

a) (m > 1) E -> (m = 1 and not *o) E -> probate -> probat cease -> ceas

b) (m > 1 and *d and *L) -> single letter

controll -> control

• The algorithm is careful not to remove a suffix when the stem is too short, the length of the stem being given by its measure, m. There is no logistic basis for this approach.

Advantages:

- Its simplicity and speed.
- Produces the best output as compared to other Stemmers.
- Less error rate.
- The Snowball Stemmer Framework designed by Porter is language independent approach to Stemming.

Limitations:

- The produced stems are not always real words.
- It has at least five steps and sixty rules and hence is time consuming.

Errors in Stemming

• There are mainly two errors in Stemming – over-stemming and under-stemming.

Over-stemming:

- Two words with different stems are stemmed to the same root.
- This is also known as a false positive.

Under-stemming:

- Two words that should be stemmed to the same root are not.
- This is also know as a false negative.

References

- Vladimír Bartík, Radek Burget
 Two-phase Categorization of Web Documents
 Brno University of Technology
- Anjali Ganesh Jivani

A Comparative Study of Stemming Algorithms The Maharaja Sayajirao University of Baroda

• William B. Frakes

Information retrieval: data structures and algorithms

References

Fuchun Peng

Context Sensitive Stemming for Web Search Sunnyvale, California

• Martin Porter

An algorithm for suffix stripping

Thank you for your attention