

Porter Stemmer

The Porter Stemmer (Porter, 1980)

- A simple rule-based algorithm for stemming
- An example of a HEURISTIC method
- Based on rules like:
 - ATIONAL -> ATE (e.g., *relational* -> *relate*)
- The algorithm consists of seven sets of rules, applied in order

The Porter Stemmer: definitions

- Definitions:
 - **CONSONANT**: a letter other than A, E, I, O, U, and Y preceded by consonant
 - **VOWEL**: any other letter
- With this definition, all words are of the form:
 $(C)(VC)^m(V)$
C=string zero or more consonants
V=string of one or more vowels
- E.g.,
 - Troubles
 - C V CVC

The Porter Stemmer: rule format

- The rules are of the form:

(condition) S1 -> S2

Where S1 and S2 are suffixes

- Conditions:

m	The measure of the stem
*S	The stem ends with S
v	The stem contains a vowel
*d	The stem ends with a double consonant
*o	The stem ends in CVC (second C not W, X, or Y)

The Porter Stemmer: Step 1

1. SSES -> SS

1. *expresses* -> *express*

2. IES -> I

1. *ponies* -> *poni*

2. *ties* -> *ti*

3. SS -> SS

1. *process* -> *process*

4. S -> ε

1. *cats* -> *cat*

The Porter Stemmer: Step 2a (past tense, progressive)

1. (m>1) EED -> EE

1. Condition verified: *agreed* -> *agree*
2. Condition not verified: *feed* -> *feed*

2. (*V*) ED -> ε

1. Condition verified: *plastered* -> *plaster*
2. Condition not verified: *bled* -> *bled*

3. (*V*) ING -> ε

- Condition verified: *motoring* -> *motor*
- Condition not verified: *sing* -> *sing*

The Porter Stemmer: Step 2b (cleanup)

- (These rules are ran if second or third rule in 2a apply)

4. **AT-> ATE**

– *conflat(ed)* -> *conflate*

5. **BL -> BLE**

– *Troubl(ing)* -> *trouble*

6. **(*d & ! (*L or *S or *Z)) -> single letter**

– Condition verified: *hopp(ing)* -> *hop*, *tann(ed)* -> *tan*

– Condition not verified: *fall(ing)* -> *fall*

7. **(m>1 & *o) -> E**

– Condition verified: *fil(ing)* -> *file*

– Condition not verified: *fail* -> *fail*

The Porter Stemmer: Steps 3 and 4

- Step 3: Y Elimination (**V**) *Y* -> *I*
 - Condition verified: *happy* -> *happi*
 - Condition not verified: *sky* -> *sky*
- Step 4: Derivational Morphology, I
 8. (*m*>0) *ATIONAL* -> *ATE*
 - *Relational* -> *relate*
 9. (*m*>0) *IZATION* -> *IZE*
 - *generalization* -> *generalize*
 10. (*m*>0) *BILITY* -> *BLE*
 - *sensibility* -> *sensible*

The Porter Stemmer: Steps 5 and 6

- Step 5: Derivational Morphology, II

- (m>0) ICATE -> IC
 - *triplicate* -> *triplic*
- (m>0) FUL -> ϵ
 - *hopeful* -> *hope*
- (m>0) NESS -> ϵ
 - *goodness* -> *good*

- Step 6: Derivational Morphology, III

- (m>0) ANCE -> ϵ
 - *allowance* -> *allow*
- (m>0) ENT -> ϵ
 - *dependent* -> *depend*
- (m>0) IVE -> ϵ
 - *effective* -> *effect*
- (m>0) IZE -> ϵ
 - *generalize* -> *general*
- (m>0) ANT -> ϵ
 - *reluctant* -> *reluct*
- (m>0) r -> ϵ
 - *computer* -> *compute*

The Porter Stemmer: Step 7 (cleanup)

- Step 7a
 - (m>1) E -> ε
 - *probate* -> *probat*
 - (m>1 & !*o) NESS -> ε
 - *goodness* -> *good*
- Step 7b
 - (m>1 & *d & *L) -> single letter
 - Condition verified: *controll* -> *control*
 - Condition not verified: *roll* -> *roll*

Examples

- *computers*
 - Step 1, Rule 4: -> *computer*
 - Step 6, Rule 4: -> *compute*
- *singing*
 - Step 2a, Rule 3: -> *sing*
- *controlling*
 - Step 2a, Rule 3: -> *controll*
 - Step 7b : -> *control*
- *generalizations*
 - Step 1, Rule 4: -> *generalization* (noun)
 - Step 4, Rule 9: -> *generalize* (verb)
 - Step 6, last rule: -> *general* (adjective)

Problems

- *elephants -> eleph*
 - Step 1, Rule 4: -> *elephant*
 - Step 6, Rule 7: -> *eleph*
- *doing - > do*
 - Step 2a, Rule 3: -> *do*

References

- The Porter Stemmer home page (with the original paper and code):
<http://www.tartarus.org/~martin/PorterStemmer/>
- Jurafsky and Martin, chapter 3.4
- The original paper: Porter, M.F., 1980, An algorithm for suffix stripping, *Program*, **14**(3):130-137.