

Candidate Roll No. _____ (In figures)

Test Exam.

Name: _____

Date: _____ 20

Examination: _____ Branch/Semester: _____

Subject: _____

Junior Supervisor's full Signature with Date	

Question No.	1	2	3	4	5	6	7	8	9	10	11	12	Total
Marks Obtained													

Stemming with Porter Stemmer Algorithm

- Suffix stripping.
- ↳ It is operation by which suffixes are removed from terms.

→ A document is represented by a vector of words, or terms, of some sort.

→ Terms with a common stem:-

- Connect → connected → connecting → connections → connect
- connected }
- connecting }
- connections }
- connect (Root word)

→ Stem is nothing but root

* Stemming In linguistic morphology and information retrieval, stemming is process of reducing terms to their "roots" before indexing.

→ Stemming programs are commonly referred to as stemming algorithms or stemmers.

→ "Stemming" suggests code or affix chopping
↳ it is language dependent-

i.e. automata(s) automatic, automation all
Defined to automat

i.e. automat is actual word, from which
all other words are generated.

∴ For Stemming one very popular algorithm
is Porter Stemmer

→ Common algorithm for Stemming English

→ The algorithm consists of 5 steps of
rules applied in order i.e. sequentially.

→ Each phase consists of set of
commands.

→ Sample convention :- of the rules in
a compound command, select one
that applies. to longest suffix

The Porter Stemmer :- definitions

Definitions :-

CONSONANT :- a letter other than A, E, I, O,
U and Y preceded by consonant.

VOWEL :- any other letter. (i.e. A, E, I, O and U)

∴ with this definition, in Porter Stemming
algorithm, all the words are of the form

(C) (VC)^m (V)

C = string of one or more consonants. (cont).

V = string of one or more vowels.

m = measure of word or word part, when
represented in this form (VC).

→ imp

Example for $M = m$ (Measure of Word)

For e.g. 1) $M = 0$; TREE, TIBBY, i.e. TR. will not have vowel consonant combination.

2) $M = 1$; TROUBLE has TRE(E), IVY.

TREE. ∵ no combination of Vc.

IVY. ∵ M = 0.

3) $M = 2$; TROUBLE has TRE(E), OATEN, ORRERY.

TROUBLE has two combinations of Vc.

∴ one combination ∵ M = 1.

4) $M = 2$; TROUBLE has PROVATE, OATEN, ORRERY.

TROUBLE has two combinations of Vc.

∴ two combinations in one word ∵ M = 2.

Rule for removing a suffix &

The rules are of the form:-

(condition) $S_1 \rightarrow S_2$ if happy cases in our form.

where condition followed by suffix S_1 maps to suffix S_2 .

where S_1 and S_2 are suffixes.

This means that if a word ends with a suffix S_1 , and the stem before S_1 satisfies the given condition then S_1 is replaced by S_2 .

→ E.g. if rule is $(m > 1)$ EMENT →

→ In this s_1 is EMENT and s_2 is NULL.

→ So this would map REPLACEMENT to REPLAC

$m = 4$ i.e. $M > 1$

satisfied.

∴ EMENT is replaced by NULL value

* The condition Patl-maj also contains One.
following $(m > 1)$ condition is true

Value	m	The measure of the stem.
*	s	The stem ends with s .
*	V	The stem contains a vowel.
*	d	The stem ends with double consonant (TT, SS).
OR		The stem ends in CVC (second c should not be W, X or Y) e.g. WIL, HOP.

→ The condition Patl-maj also contains expressions with and, or and not.

→ e.g. $(m > 1)$ and $(*s \text{ or } *t))$: tests for a stem with $m > 1$ ending in s or t .

K. J. SOMAIYA COLLEGE OF ENGINEERING

(Autonomous College Affiliated to University of Mumbai)

⑥

Candidate Roll No. _____ (In figures)

Test Exam.

Name : _____

Date : _____ 20

Examination : _____ Branch/Semester _____

Subject : _____

Junior Supervisor's full
Signature with Date

Question No.	1	2	3	4	5	6	7	8	9	10	11	12	Total
Marks Obtained													

→ The Porter & Stemmer : Step 1 a (Step 1, part 1).

i) SSES → SS

↳ effect applying Porter & Stemmer.

caresses → caress

bold → bold writing (for writing)

ii) IES → EI

↳ Ponies → Poni

ties → tie (for writing)

iii) SS → S, (for writing)

↳ Caress → caress (for writing)

iv) S → ε

↳ null (for writing)

↳ cats → cat

→ Step 1 b { with condition }

i) ($m > 0$) EED → EE

↳ (for m > 0) (for m > 0) ↳

• Condition Verified: agreed → agreed

i.e. agreed ← (and) NOT ↳

v.c. ↳

except -

↳ only for stem we have $m=1$ i.e. $m > 0$.

∴ EED is replaced by EE.

- condition not verified: ~~b~~ feed → Feed.
 ↓
 No combin' of Vc.
 ∴ no change.

ii) $(*\backslash *) ED \rightarrow \epsilon$
 \hookrightarrow_{NULL}

- condition verified: plastered → Plaster

↓
 √ $\leftarrow 2322$ ↗

→ vowel /ə/ after ED (appended by Null).
 $229807 \rightarrow 229807$

- condition not verified bled → bled
 $\hookrightarrow_{NULL} 231$ (ii)

iii) $(*\backslash *) ING \rightarrow \epsilon i t \leftarrow 2917$ ↗

- condition verified! - motoring → motory
- condition not verified: sing → sing.

3 ← 2 (vi)

* Next Step 1 cleaning (clean up)

These rules are applied if if second or third m.s.
in 1b apply).

i) AT → ATE

↳ Conflat (ed) → Conflat

ii) BLE → BLE ; BLE → BLE ↗ convert into prefix word.

↳ Troubli(ing) → Troubli(e).

iii) $(*\partial 8! (*L 08 *5 08 *2)) \rightarrow$ Single letter

i.e. stem with double consonant and final consonant -

is ends not L, S or Z then converted into single letters.

- Condition Verified: $\text{happ}(ing) \rightarrow \text{hap}$
 i.e. double consonant - then double PP replaced by single P.
 $\text{tann}(od) \rightarrow \text{tan}(son)$
 similar to first step.
 - Condition Not Verified: $\text{fail}(ing) \rightarrow \text{fail}$.
 i) $(m=1 \& *0) \rightarrow \text{fail}$
 \hookrightarrow measure = 1 and $*0$ i.e. cyc form.
 - Condition Verified: $\text{fil}(ing) \rightarrow \text{file}$
 i.e. loop \rightarrow VV loop
 - Condition Not Verified: $\text{fail} \rightarrow \text{fail}$.
 i) $(m=1 \& *0) \rightarrow \text{fail}$
 \hookrightarrow condition not satisfied.
 word \rightarrow word loop
- * The Posterior Stemmer: Step 1c and 2
- i) Step 1c: Y Elimination ($*V*$) $\rightarrow I$.
 \rightarrow loops \rightarrow switched
- Condition Verified: $\text{happy} \rightarrow \text{happi}$
 if we chop Y, then it contains a vowel.
 $\therefore Y$ is a vowel.
 - Condition Not Verified: $\text{sky} \rightarrow \text{s|ky}$ (ii)
 loop \rightarrow word loop

Step 2 Derivational Morphology, I.

- $(m>0)$ ATIONAL ($\rightarrow \text{ATE} \times \text{NOM}$)
 nouns \rightarrow relational \rightarrow relate
- $(m>0)$ IZATION \rightarrow IZE
 i.e. generalization \rightarrow generalize
- $(m>0)$ SENSIBILITY \rightarrow SENSIBLE
 \hookrightarrow Sensibiliti \rightarrow Sensible

Step 3 Derivational Morphology, II.

i) ($m > o$) ICATE \rightarrow ICAT
 \hookrightarrow triplicate \rightarrow triplic.

ii) ($m > o$) FUL \rightarrow E
 \hookrightarrow hopeful \rightarrow hope.

iii) ($m > o$) NESS \rightarrow ESK \rightarrow goodness
 \hookrightarrow goodness \rightarrow good.

Step 4 Derivational Morphology, III.

i) ($m > o$) ANCE \rightarrow E
 \hookrightarrow allowance \rightarrow allow.

ii) ($m > o$) ENT \rightarrow E
 \hookrightarrow dependent \rightarrow depend

iii) ($m > o$) IVE \rightarrow E
 \hookrightarrow effective \rightarrow effect.

Step 5 Clean Up.

i) ($m > l$) E \rightarrow E
 \hookrightarrow Probate \rightarrow Probat.

ii) ($m = 1 \vee 81! \neq o$) NESS \rightarrow ESK \rightarrow goodness
 \hookrightarrow goodness \rightarrow good.

Step 5b Clean Up

($m > 1 \neq d \neq g \neq L$) \rightarrow single letter;
 \hookrightarrow Condition Verified \rightarrow control \rightarrow controll

TSI \leftarrow TSIAS \rightarrow tsj
 \hookrightarrow Condition Not Verified \rightarrow goffl \rightarrow gojj.

This is how Porter Stemmer works on tokens and strips the roots. Woods.