

Aug 30, 07 8:28

Stemmer.java

Page 1/9

```

/*
Porter stemmer in Java. The original paper is in
5 Porter, 1980, An algorithm for suffix stripping, Program, Vol. 14,
no. 3, pp 130-137,
See also http://www.tartarus.org/~martin/PorterStemmer
10 History:
Release 1
Bug 1 (reported by Gonzalo Parra 16/10/99) fixed as marked below.
15 The words 'aed', 'eed', 'oed' leave k at 'a' for step 3, and b[k-1]
is then out outside the bounds of b.
Release 2
20 Similarly,
Bug 2 (reported by Steve Dyrdaahl 22/2/00) fixed as marked below.
'ion' by itself leaves j = -1 in the test for 'ion' in step 5, and
b[j] is then outside the bounds of b.
25 Release 3
Considerably revised 4/9/00 in the light of many helpful suggestions
from Brian Goetz of Quiotix Corporation (brian@quiotix.com).
30 Release 4
*/
35 import java.io.*;
/**
 * Stemmer, implementing the Porter Stemming Algorithm
 *
 * The Stemmer class transforms a word into its root form. The input word can be
 * provided a character at time (by calling add()), or at once by calling one of
 * the various stem(something) methods.
 */
45 class Stemmer {
    private char[] b;
    private int i, /* offset into b */
        i_end, /* offset to end of stemmed word */
        j, k;
    private static final int INC = 50;
    /* unit of size whereby b is increased */
    public Stemmer() {
55         b = new char[INC];
        i = 0;
        i_end = 0;
    }
    /**
60     * Add a character to the word being stemmed. When you are finished adding
     * characters, you can call stem(void) to stem the word.
     */
    public void add(char ch) {
65         if (i == b.length) {
            char[] new_b = new char[i + INC];
            for (int c = 0; c < i; c++)
                new_b[c] = b[c];
70             b = new_b;
        }
        b[i++] = ch;
    }
}

```

Thursday August 30, 2007

Stemmer.java

Aug 30, 07 8:28

Stemmer.java

Page 2/9

```

75 /**
 * Adds wLen characters to the word being stemmed contained in a portion of
 * a char[] array. This is like repeated calls of add(char ch), but faster.
 */
    public void add(char[] w, int wLen) {
80         if (i + wLen >= b.length) {
            char[] new_b = new char[i + wLen + INC];
            for (int c = 0; c < i; c++)
                new_b[c] = b[c];
85             b = new_b;
        }
        for (int c = 0; c < wLen; c++)
            b[i++] = w[c];
    }
90 /**
 * After a word has been stemmed, it can be retrieved by toString(), or a
 * reference to the internal buffer can be retrieved by getResultBuffer and
 * getResultLength (which is generally more efficient.)
 */
95 public String toString() {
    return new String(b, 0, i_end);
}
/**
100 * Returns the length of the word resulting from the stemming process.
 */
public int getResultLength() {
    return i_end;
}
105 /**
 * Returns a reference to a character buffer containing the results of the
 * stemming process. You also need to consult getResultLength() to determine
 * the length of the result.
 */
110 public char[] getResultBuffer() {
    return b;
}
115 /* cons(i) is true <=> b[i] is a consonant. */
    private final boolean cons(int i) {
        switch (b[i]) {
120             case 'a':
            case 'e':
            case 'i':
            case 'o':
            case 'u':
                return false;
125             case 'y':
                return (i == 0) ? true : !cons(i - 1);
            default:
                return true;
        }
130     }
    /**
 * m() measures the number of consonant sequences between 0 and j. if c is a
 * consonant sequence and v a vowel sequence, and <..> indicates arbitrary
 * presence,
135     * <c><v> gives 0 <c><vc> gives 1 <c><vcvc> gives 2 <c><vcvcvc> gives
     * 3 ....
     */
140     private final int m() {
        int n = 0;
        int i = 0;
        while (true) {
145             if (i > j)
                return n;

```

1/5

Aug 30, 07 8:28

Stemmer.java

Page 3/9

```

        if (!cons(i))
            break;
        i++;
150     }
        i++;
        while (true) {
            while (true) {
                if (i > j)
155                 return n;
                if (cons(i))
                    break;
                i++;
            }
            i++;
            n++;
            while (true) {
                if (i > j)
160                 return n;
                if (!cons(i))
165                 break;
                i++;
            }
            i++;
170     }
}

/* vowelinstem() is true <=> 0,...j contains a vowel */
175 private final boolean vowelinstem() {
    int i;
    for (i = 0; i <= j; i++)
        if (!cons(i))
            return true;
180     return false;
}

/* doublec(j) is true <=> j,(j-1) contain a double consonant. */
185 private final boolean doublec(int j) {
    if (j < 1)
        return false;
    if (b[j] != b[j - 1])
        return false;
190     return cons(j);
}

/*
 * cvc(i) is true <=> i-2,i-1,i has the form consonant - vowel - consonant
 * and also if the second c is not w,x or y. this is used when trying to
 * restore an e at the end of a short word. e.g.
 * cav(e), lov(e), hop(e), crim(e), but snow, box, tray.
200 */

private final boolean cvc(int i) {
    if (i < 2 || !cons(i) || cons(i - 1) || !cons(i - 2))
        return false;
205     {
        int ch = b[i];
        if (ch == 'w' || ch == 'x' || ch == 'y')
            return false;
    }
    return true;
210 }

private final boolean ends(String s) {
    int l = s.length();
    int o = k - l + 1;
    if (o < 0)
        return false;
215     for (int i = 0; i < l; i++)
        if (b[o + i] != s.charAt(i))

```

Aug 30, 07 8:28

Stemmer.java

Page 4/9

```

        return false;
        j = k - 1;
        return true;
    }
225     /*
     * setto(s) sets (j+1),...k to the characters in the string s, readjusting
     * k.
     */

230     private final void setto(String s) {
        int l = s.length();
        int o = j + 1;
        for (int i = 0; i < l; i++)
            b[o + i] = s.charAt(i);
235     }

    /* r(s) is used further down. */

240     private final void r(String s) {
        if (m() > 0)
            setto(s);
    }

245     /*
     * step1() gets rid of plurals and -ed or -ing. e.g.
     * caresses -> caress ponies -> poni ties -> ti caress -> caress cats -> cat
     *
     * feed -> feed agreed -> agree disabled -> disable
     *
     * matting -> mat mating -> mate meeting -> meet milling -> mill messing ->
     * mess
     *
     * meetings -> meet
     */

255     private final void step1() {
        if (b[k] == 's') {
            if (ends("sses"))
                k -= 2;
            else if (ends("ies"))
                setto("i");
260             else if (b[k - 1] != 's')
                k--;
        }
        if (ends("eed")) {
            if (m() > 0)
                k--;
270         } else if ((ends("ed") || ends("ing")) && vowelinstem()) {
            k = j;
            if (ends("at"))
                setto("ate");
            else if (ends("bl"))
                setto("ble");
            else if (ends("iz"))
                setto("ize");
            else if (doublec(k)) {
280                 k--;
                {
                    int ch = b[k];
                    if (ch == 'l' || ch == 's' || ch == 'z')
                        k++;
                }
            } else if (m() == 1 && cvc(k))
                setto("e");
        }
    }

290     /* step2() turns terminal y to i when there is another vowel in the stem. */

```

Aug 30, 07 8:28

Stemmer.java

Page 5/9

```

private final void step2() {
    if (ends("y") && vowelinstem())
295         b[k] = 'i';
}
/*
 * step3() maps double suffices to single ones. so -ization (= -ize plus
 * -ation) maps to -ize etc. note that the string before the suffix must
 * give m() > 0.
 */
private final void step3() {
305     if (k == 0)
        return; /* For Bug 1 */
    switch (b[k - 1]) {
    case 'a':
310         if (ends("ational")) {
            r("ate");
            break;
        }
        if (ends("tional")) {
315             r("tion");
            break;
        }
        break;
    case 'c':
320         if (ends("enci")) {
            r("ence");
            break;
        }
        if (ends("anci")) {
325             r("ance");
            break;
        }
        break;
    case 'e':
330         if (ends("izer")) {
            r("ize");
            break;
        }
        break;
    case 'l':
335         if (ends("bli")) {
            r("ble");
            break;
        }
        if (ends("alli")) {
340             r("al");
            break;
        }
        if (ends("entli")) {
345             r("ent");
            break;
        }
        if (ends("eli")) {
350             r("e");
            break;
        }
        if (ends("ousli")) {
355             r("ous");
            break;
        }
        break;
    case 'o':
360         if (ends("ization")) {
            r("ize");
            break;
        }
        if (ends("ation")) {
365             r("ate");
            break;
        }
        if (ends("ator")) {

```

Aug 30, 07 8:28

Stemmer.java

Page 6/9

```

        r("ate");
        break;
    }
    break;
370     case 's':
        if (ends("alism")) {
            r("al");
            break;
        }
375         if (ends("iveness")) {
            r("ive");
            break;
        }
        if (ends("fulness")) {
380             r("ful");
            break;
        }
        if (ends("ousness")) {
385             r("ous");
            break;
        }
        break;
    case 't':
390         if (ends("aliti")) {
            r("al");
            break;
        }
        if (ends("iviti")) {
395             r("ive");
            break;
        }
        if (ends("biliti")) {
400             r("ble");
            break;
        }
        break;
    case 'g':
405         if (ends("logi")) {
            r("log");
            break;
        }
    }
}
410 /* step4() deals with -ic-, -full, -ness etc. similar strategy to step3. */
private final void step4() {
    switch (b[k]) {
    case 'e':
415         if (ends("icate")) {
            r("ic");
            break;
        }
        if (ends("ative")) {
420             r("");
            break;
        }
        if (ends("alize")) {
425             r("al");
            break;
        }
        break;
    case 'i':
430         if (ends("iciti")) {
            r("ic");
            break;
        }
        break;
    case 'l':
435         if (ends("ical")) {
            r("ic");
            break;
        }
    }
}

```

Aug 30, 07 8:28

Stemmer.java

Page 7/9

```

440     if (ends("ful")) {
        r("");
        break;
    }
    break;
445 case 's':
    if (ends("ness")) {
        r("");
        break;
    }
450 }
}

/* step5() takes off -ant, -ence etc., in context <c>vcvc<v>. */

455 private final void step5() {
    if (k == 0)
        return; /* for Bug 1 */
    switch (b[k - 1]) {
460 case 'a':
        if (ends("al"))
            break;
        return;
    case 'c':
465     if (ends("ance"))
            break;
        if (ends("ence"))
            break;
        return;
    case 'e':
470     if (ends("er"))
            break;
        return;
    case 'i':
475     if (ends("ic"))
            break;
        return;
    case 'l':
480     if (ends("able"))
            break;
        if (ends("ible"))
            break;
        return;
    case 'n':
485     if (ends("ant"))
            break;
        if (ends("ement"))
            break;
        if (ends("ment"))
            break;
490     /* element etc. not stripped before the m */
        if (ends("ent"))
            break;
        return;
    case 'o':
495     if (ends("ion") && j >= 0 && (b[j] == 's' || b[j] == 't'))
            break;
        /* j >= 0 fixes Bug 2 */
        if (ends("ou"))
            break;
500     return;
        /* takes care of -ous */
    case 's':
505     if (ends("ism"))
            break;
        return;
    case 't':
        if (ends("ate"))
            break;
        if (ends("iti"))
            break;
510     return;

```

Aug 30, 07 8:28

Stemmer.java

Page 8/9

```

    case 'u':
        if (ends("ous"))
            break;
515     return;
    case 'v':
        if (ends("ive"))
            break;
        return;
520     case 'z':
        if (ends("ize"))
            break;
        return;
    default:
525     return;
}
if (m() > 1)
    k = j;
}

530 /* step6() removes a final -e if m() > 1. */

private final void step6() {
    j = k;
535     if (b[k] == 'e') {
        int a = m();
        if (a > 1 || a == 1 && !cvc(k - 1))
            k--;
    }
540     if (b[k] == 'l' && doublec(k) && m() > 1)
        k--;
}

/** Stem the word placed into the Stemmer buffer through calls to add().
 * Returns true if the stemming process resulted in a word different
 * from the input. You can retrieve the result with
 * getResultLength()/getResultBuffer() or toString().
 */
550 public void stem() {
    k = i - 1;
    if (k > 1) {
        step1();
        step2();
        step3();
555     step4();
        step5();
        step6();
    }
    i_end = k + 1;
560     i = 0;
}

/** Test program for demonstrating the Stemmer. It reads text from a
 * a list of files, stems each word, and writes the result to standard
 * output. Note that the word stemmed is expected to be in lower case:
 * forcing lower case must be done outside the Stemmer class.
 * Usage: Stemmer file-name file-name ...
 */
565 public static void main(String[] args) {
    char[] w = new char[501];
    Stemmer s = new Stemmer();
    for (int i = 0; i < args.length; i++)
        try {
570             FileInputStream in = new FileInputStream(args[i]);

            try {
                while (true)
                    {
580                         int ch = in.read();
                            if (Character.isLetter((char) ch)) {
                                int j = 0;
                                while (true) {
                                    ch = Character.toLowerCase((char) ch);

```

Aug 30, 07 8:28

Stemmer.java

Page 9/9

```

585         w[j] = (char) ch;
        if (j < 500)
            j++;
        ch = in.read();
590         if (!Character.isLetter((char) ch)) {
            /* to test add(char ch) */
            for (int c = 0; c < j; c++)
                s.add(w[c]);

            /* or, to test add(char[] w, int j) */
            /* s.add(w, j); */

595         s.stem();
            {
                String u;

600                 /* and now, to test toString() : */
                u = s.toString();

                /* to test getResultBuffer(), getResultL
605         length() : */
                /* u = new String(s.getResultBuffer(), 0
                , s.getResultLength()); */

                System.out.print(u);
            }
610         }
        }
        if (ch < 0)
            break;
615         System.out.print((char) ch);
    }
    } catch (IOException e) {
        System.out.println("error reading " + args[i]);
        break;
620    }
    } catch (FileNotFoundException e) {
        System.out.println("file " + args[i] + " not found");
        break;
625    }
}
}

```