

STYLOCHRONOMETRY WITH SUBSTRINGS
Or: A Poet Young and Old

Richard S. Forsyth.

Department of Computing & Information Systems,
Faculty of Science, Design & Technology,
University of Luton,
Luton LU1 3JU.
U.K.

[
Contact:
forsyth_rich@yahoo.co.uk

Cite as:

Forsyth, R.S. (1999). Stylochronometry with substrings, or: a poet young and old. *Literary & Linguistic Computing*, 14(4), 1-11.

]

STYLOCHRONOMETRY WITH SUBSTRINGS Or: A Poet Young and Old

"Laughter not time destroyed my voice
And put that crack in it"
-- *A Man Young and Old*, W.B. Yeats (1926).

Abstract

Assigning a date to a text is an important task in stylometry. Most previous researchers, however, have worked on intractable problems, where a true chronology will never be known with certainty, such as the works of Plato, Shakespeare or Marlowe. It is argued here that stylochronometric methods should be extensively tested on unproblematic texts before being used in disputed cases. As part of such testing, the present study applies a novel technique, Monte-Carlo Feature-Finding (Forsyth, 1997), to the verse of W.B. Yeats, where the dating is relatively well documented. Yeats insisted that his language changed as he grew older, and most readers would concur; yet scholars have not reached agreement on the nature of this linguistic change (Jaynes, 1980).

A quasi-random search algorithm was used to find marker substrings in 142 poems of Yeats. To test their distinctiveness, four trials were performed: (1) assignment of 10 poems absent from the training sample to their correct period; (2) detecting differences in two poems written by Yeats in his twenties and revised when he was sixty; (3) constructing a regression formula; (4) classifying two prose extracts written 46 years apart.

Assigning short poems (median length 114 words) to their correct chronological period is a non-trivial task. Nevertheless, counting of distinctive substrings gave the right assignment in 9 out of 10 unseen cases. Moreover, these substring frequencies were sensitive enough to detect authorial revision in two early poems revised by Yeats many years after he originally wrote them, and robust enough to classify a pair of short prose extracts correctly; as well as accounting for 71% of the variance when used in a regression to predict the year in which 13 poems, absent from the training sample, were composed.

These results suggest that short substrings found by a Monte-Carlo process warrant further investigation as stylistic indicators.

Keywords: Monte-Carlo Methods, Stylistics, Stylochronometry, Yeats.

1. Introduction

Methods of assigning dates to texts have long been of interest to literary scholars. In the words of Graham Martin:

"in one way or another, the study of even a single poet moves necessarily towards a process of comparison and contrast between poem and poem and then between poet and poet, and because poets write and publish at particular moments in time, that means, at some point, chronology."
(Martin, 1975, p. 5)

Apart from authorship attribution, date assignment (stylochronometry) is the most commonly reported task in the stylometric literature. Yardi (1946), Brainerd (1980), Ule (1982), Foster (1989) and Temple (1996) are just some of the researchers who have addressed this problem.

A variety of different linguistic indicators have been used for this task, often based on prosody or scansion (such as the proportion of 'feminine' line endings), and many such indicators are generally used in combination, as part of a multivariate analysis. Nevertheless, Foster reported good results with a single variable:

"the frequency of *most* (excluding its use as a substantive) is one of the best indicators for dating Shakespeare's work, and certainly one of the simplest." (Foster, 1989, p. 109)

Methodologically, however, the five studies cited above suffer from a common problem: they deal with cases where a true chronology will never be known with certainty -- the plays of Shakespeare, the writings of Christopher Marlowe and the Platonic corpus. The present study, in contrast, follows the practice of Frischer (1991) in his attempt to assign a date to Horace's *Ars Poetica*. Frischer first conducted a systematic search through the poet's securely dated works for what he termed chronometers -- habits that change regularly over time -- and then applied these chronometers (such as the usage rates of 'ad' or 'sed') to estimate the unknown date of the *Ars Poetica*. Statistical methods were used to filter out from the potential indicators those whose pattern of temporal variation could be explained as chance variation.

The present investigation follows this general approach, and thus differs from much previous research in this field, by concentrating on a test case -- the verse of William Butler Yeats -- where the dating is comparatively well attested. Additionally, it extends Frischer's paradigm by further automating the search for potential chronometers. Specifically, it applies a novel technique called Monte-Carlo Feature-Finding (Forsyth & Holmes, 1996) to test its suitability in this area. This method seeks short substrings and accordingly is not limited to discovering **words** as chronometers. It can also find subwords (e.g. affixes such as 'ing' or 'tion') as well as collocations (such as 'in the' or 'ed by').

Dating Yeats makes an interesting initial test because over the course of a long poetic career his style changed noticeably; yet scholars do not agree on the nature of that change. In this context, it is pertinent to quote Jaynes (1980), who also studied the

evolution of Yeats's poetic style:

"Yeats's syntactic style is quite stable over his career." (Jaynes, 1980, p. 13)

as well as Martindale (1990), who investigated historical trends in poetic language in both English and French over several centuries:

"the content of most poets' verse does not change massively across the course of a lifetime".

However, Yeats himself firmly believed that his syntax and diction evolved as he grew older (a view shared by many of his readers) as the following quotations testify.

"My own verse has more and more adopted -- seemingly without any will of mine -- the syntax and vocabulary of common personal speech." (Yeats, 1926; in Jaynes, 1980, p. 11)

"It was a long time before I had made a language to my liking". (Yeats, 1937; in Jeffares, 1964, p. 265)

2. Method

The main aim of the present study was to find out whether an automated method of characterizing textual differences which has shown promise in other areas (Forsyth, 1997) could be used to analyze the development in a particular poet's use of language. This method, Monte-Carlo Feature-Finding, is simply a random search for substrings that exist in the training data. It is implemented by a program called TEFF (Text Extending Feature Finder). This finds markers merely by searching through a given set of training texts.

TEFF randomly picks many substrings (4096 in the experiment reported here) from the combined training text. The length of each substring is also a random number from 1 to 8. All distinct substrings thus found are saved and then ranked according to their distinctiveness, i.e. according to a measure of their differential rate of occurrence in the different text categories for the problem concerned. Chi-squared is used to measure distinctiveness:

$$\sum_j \frac{(O_j - E_j)^2}{E_j}$$

where O_j is the observed frequency in a given text category and E_j is the expected frequency, assuming an equal rate of occurrence in each category.

Earlier trials of Monte-Carlo Feature-Finding (Forsyth, 1995) indicated that the basic procedure, as described above, tends to generate substrings that are fragmented at what seem linguistically inappropriate boundary points, even when they prove effective as discriminators. This weakness has since been rectified, by a 'string stretching' procedure. For technical details see Forsyth & Holmes (1996) or Forsyth (1997).

Another improvement to the basic TEFF procedure was employed in the experiment reported here. The texts were divided into blocks of approximately 1030 characters (in fact at the first white-space character after 1024 bytes). Then substrings were counted within blocks and the square-root of each within-block count summed to give an adjusted total frequency for each substring. The practical effect of this adjustment is to downgrade the importance of indicators that occur in localized bursts relative to those whose usage is pervasive throughout a text. (A fuller explanation is given in Forsyth (1997).)

3. Materials

A random sample of 142 poems by W.B. Yeats was used as training data. This sample was divided into two portions:

"Younger Yeats" (YY), 72 poems, 18,360 words;
 "Older Yeats" (OY), 70 poems, 18,668 words.

The dividing date was 1915. Thus the YY sample was written in 1914 or before and the OY sample from 1916 onwards¹.

The TEFF program was used to find substrings which distinguished these two classes. Then the efficacy of these substrings as distinctive markers was tested on three types of unseen material: (1) 10 other poems, absent from the training sample; (2) two early poems that Yeats later heavily revised; and (3) two prose extracts. Details of these tests are given in the next section.

4. Results

The following listing (Table 1) shows the 65 most distinctive substrings found by the TEFF program when given samples of poetry written by William Butler Yeats before 1915 ('Younger Yeats') or after 1915 ('Older Yeats').

Table 1 -- Substrings Found by TEFF.

TEFF output; date: 10/25/97 21:27:58

Rank	Substring	Chi-score		
1	`what	35.2977123	30.	100.
2	` can	34.5481048	21.	82.
3	` can	32.1377981	13.	66.
-4	`hat	29.8459659	132.	245.
5	`hat	28.7157873	126.	235.
6	` whi	25.8148745	67.	21.
7	`s, an	25.3604502	63.	19.
8	` sea	23.8708713	52.	13.
9	` that	23.0786873	114.	206.
10	`?	22.3848552	30.	83.
11	` with	22.242135	139.	74.
12	` int	21.9169178	12.	51.
13	` . ii	21.4710501	0.	23.
14	` stars	20.9510074	26.	2.

¹ This division is not haphazard: 1914 is the date of the outbreak of World War I and 1916 the date of the Irish Easter Rising, both events that profoundly affected Yeats's beliefs about what poetry should aim to achieve.

15	` that	20.8584369	114.	200.
16	`ith	20.3048191	134.	72.
17	`s of	20.0069369	105.	52.
18	`e that	19.4328782	18.	58.
19	`though	18.6119958	30.	77.
20	` you	18.4021636	120.	65.
21	`ou	18.338578	70.	29.
22	`e that	17.6827819	18.	56.
-23	`at	17.6339928	223.	332.
24	`, and	17.5216778	173.	108.
-25	`ith	17.5057402	143.	84.
-26	` tha	17.4767001	127.	209.
27	`ping	17.3800959	34.	7.
28	` you	17.1449342	68.	29.
29	`your	16.933698	52.	19.
30	` we	16.8084892	141.	83.
-31	`, an	16.6399743	176.	112.
32	` your	16.5048223	49.	17.
33	`woo	16.2479445	32.	7.
34	`ck	15.9644919	51.	103.
35	`low	15.935057	70.	32.
36	`wed	15.7755436	23.	3.
37	` the w	15.7503404	103.	56.
38	`w the	15.5959647	31.	7.
39	`murmur	15.5603536	21.	2.
40	` of o	15.5309625	27.	5.
41	` dee	15.3281468	28.	5.
42	` wa	15.2033782	141.	86.
43	`ec	15.0040229	42.	89.
44	`its	14.9342107	16.	49.
45	`io	14.8929466	46.	94.
46	`hat c	14.7247314	11.	39.
47	` fl	14.6750908	85.	44.
48	`ic	14.641333	54.	105.
49	`. an	14.4836672	39.	12.
-50	`th	14.3405652	155.	99.
51	`mag	14.2989514	9.	36.
-52	`wa	14.263478	168.	110.
53	`tion	14.2141891	7.	32.
54	`ountain	14.0862749	0.	16.
-55	`s,	13.5564112	160.	105.
-56	`he w	13.4639408	110.	65.
57	` over th	13.2886104	16.	1.
58	`ering	13.149059	61.	28.
-59	`t	13.1179103	314.	426.
60	`t can	12.8095651	2.	19.
61	`mat	12.8064674	3.	22.
62	`e white	12.5426343	13.	0.
63	`k and	12.519069	2.	19.
64	`ittle	12.4591915	24.	5.
65	`the fenian	12.3657293	13.	0.

Roughly speaking, a string is a Younger-Yeats marker if the first of the last 2 frequency counts in the preceeding table is higher than the second; so, for example, `murmur', at rank 39, is a Younger-Yeats marker (n=21 versus n=2).

A minus-sign in front of a rank number indicates a marker which is a proper substring of another higher up in the table, and which will therefore not be saved on file. This removes some of the redundancy from the list, but not all of it. So a further program was written to filter out substrings that duplicate the effect of those higher in the list. This was done by `destructive' counting: in TEFF, substrings are counted independently, but in the post-filtering program they are read in the order output by TEFF (an order of distinctiveness) and processed in that order. As a string is counted,

it is removed from the text block currently being processed. So, for example, if `what' has already been counted, occurrences of `hat' will only be found that are not embedded within `what' (as in `that' or `hatred' but not `whatever'). The result of running this latter program is shown in Table 2.

Table 2 – Marker Substrings.

```
filters.spt date: 10/25/97 22:07:45
1 C:\BM95\WY.Yx 96547 bytes.
2 C:\BM95\WY.Ox 100322 bytes.
proportion in class 1 = 0.4904112
proportion in class 2 = 0.5095863
Grams kept = 88
```

Rank	Substring	Chi-score		
1	`what	35.1125492	30.	100.
2	` can	34.3848949	21.	82.
3	`s, an	25.4882859	63.	19.
4	` whi	25.4359635	67.	21.
5	` with	22.3028559	139.	74.
6	`?	21.9694592	30.	83.
7	` sea	21.9076975	49.	13.
8	` int	21.814545	12.	51.
9	`. ii	21.4093664	0.	23.
10	` stars	21.0202576	26.	2.
11	`ck	20.7388396	35.	87.
12	` we	20.6404509	139.	76.
13	`hat	20.6188054	113.	200.
14	`s of	19.2516067	104.	52.
15	`though	18.5468022	30.	77.
16	` you	18.1662318	115.	62.
17	` that	17.7369649	3.	27.
18	`ping	16.5610752	33.	7.
19	`woo	16.3187661	32.	7.
20	` dee	16.2970505	27.	4.

It should be noted that these are strings, not words, so item 16 ` you' (a YY marker) would cover words such as `young' `younger' `youth' and `your' as well as the word `you' itself.

The presence of `?' on this list suggests that this method is potentially capable of detecting some kinds of syntactic or semiotic change as well as changes in vocabulary. However, a question-mark may not consistently signal the same syntactic structure (still less the same kind of speech act) over the course of an author's career; so the interpretation of such a change in preference still requires scrutiny of the texts themselves.

Only the top twenty markers have been shown in Table 2, to save space; and only these will be used in the analysis that follows. Note that the sole element of human judgement involved in choosing these markers was deciding how many to use, 20 being a convenient number.

5. Some Marks of Time

As an illustration, Table 3 shows the results of counting the occurrences of eleven Younger-Yeats and nine Older-Yeats marker substrings in a pair of poems written 50 years apart.

Table 3 -- Frequencies of Substrings in 2 Short Poems.

Marker Substrings:	Salley Gardens 1888 [98 words]	Politics 1938 [80 words]
Younger-Yeats:		
`s, an`	0	0
` whi`	0	0
` with`	2	0
` sea`	0	0
` stars`	0	0
` we`	1	1
`s of`	0	0
` you`	2	1
`ping`	0	0
`woo`	0	0
` dee`	0	0
Total =	5	2
Older-Yeats:		
`what`	0	2
` can`	0	1
`?`	0	1
` int`	0	0
` . ii`	0	0
`ck`	0	0
`hat`	0	6
`though`	0	1
` that`	0	4
Total =	0	15

This shows a clear preponderance of `younger' markers in the earlier poem and an even clearer preponderance of `older' markers in the later poem. If either of these texts had just been rediscovered, we could with reasonable confidence allocate *Down by the Salley Gardens* to Yeats's early career and *Politics* to his later years.

Perhaps it is worth noting here that stylometers have generally found chronology a trickier subject than authorship attribution (Forsyth, 1995), and that they have very rarely dared to categorize text segments as short as these 2 poems -- a notable exception being Simonton (1990), who analyzed, among other things, word usage in the final couplets of Shakespeare's sonnets, averaging 17.6 words in length.

However, both these poems were present in the training files. Thus Table 3 merely illustrates this method.

5.1 Unseen Trial

As a genuine test, 10 more poems were chosen from Yeat's Collected Poems, at random, five written before 1915 and five afterwards -- with the proviso that they were not already present in the training sample. These poems, and their dates of composition, are listed in Table 4, along with counts of the YY and OY substrings found in each.

Table 4 -- Substring Counts in 10 Unseen Poems.

Poem ↓	Total Count of Younger-Yeats Markers (YY)	Total Count of Older-Yeats Markers (OY)
A Faery Song (1891) [104 words]	10	1
The Lover Tells of the Rose in His Heart (1892) [114 words]	10	4
The Hosting of the Sidhe (1893) [124 words]	6	2
The Host of the Air (1893) [310 words]	23	4
To Some I Have Talked with by the Fire (1895) [139 words]	9	0
A Woman Young and Old -- Parting (1926) [79 words]	4	6
In Memory of Eva Gore-Booth and Con Markiewicz (1927) [190 words]	8	6
Quarrel in Old Age (1931) [72 words]	1	11
Parnell's Funeral, Part I (1933) [247 words]	4	33
A Model for the Laureate (1937) [142 words]	4	18

In nine of the 10 poems the count is higher in the appropriate age category. The probability of 9 or more correct binary choices from 10, under a Null Hypothesis that there is an even chance of being right or wrong, is $11/1024$ ($p=0.0107$), suggesting that short substrings can indeed be useful in this sort of problem.

5.2 A Youthful Yeatsian Index

Of course ideally in stylochronometry we would want not just to classify texts as early versus late, but to assign an estimated date to each text. With this object in mind, a 'youthful Yeatsian index' (YYIX) was defined as follows

$$YYIX = (YY - OY) / (YY + OY)$$

where YY is the number of younger-Yeats markers found and OY the number of older-Yeats markers. (This time only 19 substrings were used -- omitting 'hat' from the list shown in Table 3 on the grounds that both 'what' and 'that' were already present - - to avoid any suggestion of double-counting.) In addition, three more unseen poems were added to the 10 listed in Table 4, to give a more balanced distribution across Yeats's career, including his middle period. These were *The Ragged Wood* (1904, 105 words), *All Things can Tempt me* (1908, 92 words) and *The Scholars* (1915, 73 words).

Figure 1, below, shows a plot of YYIX against date of composition. The correlation of YYIX with date is $r = -0.844$ which is very highly significant ($p < 0.001$). A regression using YYIX to predict date accounts for over 71% of the variance.

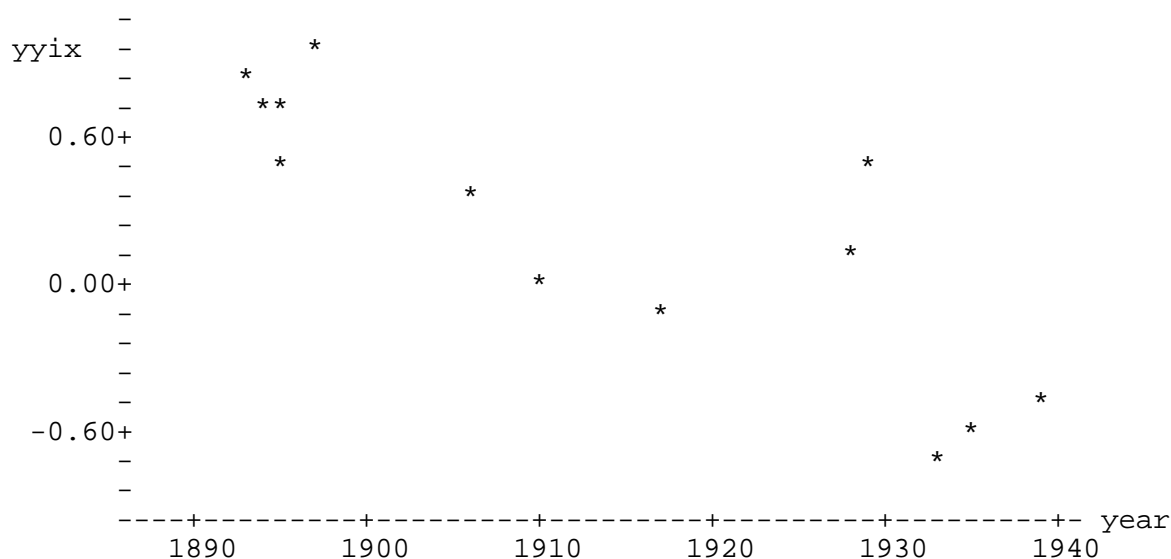


Figure 1 -- Plot of Youthful Yeatsian Index against Date.

5.3 Visions and Revisions

Another test of this approach involved looking at poems revised by Yeats in his later years. In the edition used here (Albright, 1990) there are just two such poems where the text of both versions is given in full. These are *The Lamentation of the Old Pensioner* (original 1890, revised 1925) and *The Sorrow of Love* (original 1891, revised 1924).

Table 5 gives the frequencies of YY and OY substrings in both versions of both these poems.

Table 5 -- Substring Frequencies in Revised Poems.

	YY markers	OY markers
The Lamentation ... 1890 version	2	2
1925 version [95 words]	0	11
The Sorrow of Love 1891 version	13	0
1924 version [96 words]	2	6

In each case Yeats's process of revision increases the number of `older-Yeats' markers and decreases the number of `younger-Yeats' markers.

5.4 Prose Trial

Although the training files used by TEFF to find distinctive substrings were entirely composed of poetry, it was thought worthwhile to look briefly at the frequencies of these marker substrings in prose passages as well, to gain an initial idea of the robustness of this approach when assumptions about genre are violated. Accordingly, two short extracts of prose were taken from the first and last essay by Yeats in the collection of Jeffares (1964), namely the first 40 lines of *The Irish National Literary Society* (446 words, dated 1892) and the first 44 lines of *Ireland after the Revolution* (435 words, dated 1938).

In the earlier extract, there were 21 YY markers and 11 OY markers. The later passage contained 11 YY markers and 36 OY markers. Of the 11 YY substrings, 7 were more frequent in the earlier piece than the later (with 3 tied in frequency). Of the 9 OY substrings, four were more frequent in the later piece than the earlier (with 4 tied). To put this another way: only 2 of the 20 markers pointed in the `wrong' direction.

6. Discussion

Counting `badges' in this manner is a rather unsophisticated method of text classification, so the performance of the marker substrings found by TEFF on the trials described above is quite impressive, especially bearing in mind the unreliability of most previous stylometric techniques on samples as small as those analyzed in this study -- as witnessed (for example) by the following quotation.

"Even using 500 word samples we should anticipate a great deal of unevenness" (Ledger & Merriam, 1994).

Weaknesses still remain. Firstly, the presence of `hat ' as well as ` that' in the list of OY markers suggests that the post-filtering program still needs improvement. There is clear overlap between these two substrings, which introduces an undesirable element of double-counting. Secondly, and perhaps more important, interpretation of markers like `though' is problematic. The question arises: does a group of words such as

`though', `although', `thought' and `thoughtful' constitute a *natural kind*? And if not, are we justified in relying on such a heterogenous grouping? The answer must depend, in large part, on what we want to do with the texts under scrutiny. If insight is our aim, then a KWIC index based on the substrings in question should enlighten us about just what verbal habits are being detected. If accuracy of estimation is our prime concern, the standard statistical technique of cross-validation (Weiss & Kulikowski, 1991) should protect us from jumping to spurious conclusions.

7. Conclusion

To conclude: assigning short poems by W.B. Yeats to their correct chronological period is a non-trivial task. Nevertheless, a simple count of distinctive substrings found by the TEFF program led to the right assignment in 9 out of 10 unseen cases. Moreover, an easily calculated numerical index (YYIX) contrasting younger-Yeats with older-Yeats markers showed a highly significant correlation ($r = -0.844$) with the date of composition of 13 poems (median length = 114 words) in a genuine out-of-sample test. In addition, these substring frequencies were sensitive enough to detect authorial revision in two early poems revised by Yeats many years after he originally wrote them, and robust enough to classify a pair of short prose extracts correctly as well.

The performance in this pilot study of short substrings found by a Monte-Carlo process suggests that such strings warrant further investigation as stylistic indicators.

Further work along these lines is planned. Firstly, it is necessary to check whether this empiricist approach works well with other genres, e.g. prose, with other authors, and preferably also in other languages. Assuming these further trials confirm that this approach does work in practice, it will then become desirable to explore ways of linking the low-level markers found by the Monte-Carlo process with higher-level linguistic constructs, in order to make the results found easier to interpret.

References

- Albright, D. (1990) ed. *The Poems: W.B. Yeats*. Everyman/Dent, London.
- Brainerd, B. (1980). The Chronology of Shakespeare's Plays: a Statistical Study. *Computers & the Humanities*, 14, 221-230.
- Forsyth, R.S. (1995). *Stylistic Structures: a Computational Approach to Text Classification*. Unpublished PhD thesis, Faculty of Science, University of Nottingham.
- Forsyth, R.S. (1997). Deriving Document Descriptors from Data. In: L. Dorfman et al. (eds.) *Emotion, Creativity and Art*. Perm, Russia.
- Forsyth, R.S. & Holmes, D.I. (1996). Feature-Finding for Text Classification. *Literary & Linguistic Computing*, 11(4), 163-174.
- Foster, D.W. (1989). *Elegy by W.S. A Study in Attribution*. University of Delaware Press, Newark.
- Frischer, B. (1991). *Shifting Paradigms: New Approaches to Horace's Ars Poetica*. Scholars Press, Atlanta GA.
- Jaynes, J.T. (1980). A Search for Trends in the Poetic Style of W.B. Yeats. *ALLC Journal*, 1, 11-18.
- Jeffares, A.N. (1964) ed. *Yeats: Selected Criticism*. Macmillan & Co. London.
- Ledger, G.R. & Merriam, T.V.N. (1994). *Shakespeare, Fletcher, and the Two Noble*

- Kinsmen. *Literary & Linguistic Computing*, 9(4), 235-248.
- Martin, G. (1975). *English Poetry in 1912*. Open University Press, Milton Keynes.
- Martindale, C. (1990). *The Clockwork Muse*. Basic Books, New York.
- Simonton, D.K. (1990). Lexical Choices and Aesthetic Success: a Computer Content Analysis of 154 Shakespeare Sonnets. *Computers & the Humanities*, 24, 251-264.
- Temple, J.T. (1996). A Multivariate Synthesis of Published Platonic Stylometric Data. *Literary & Linguistic Computing*, 11(2), 67-75.
- Ule, L. (1982). Recent Progress in Computer Methods of Authorship Determination. *ALLC Bulletin*, 10(3), 73-89.
- Weiss, S.M. & Kulikowski, C.A. (1991). *Computer Systems that Learn*. Morgan Kaufmann, San Mateo, CA.
- Yardi, M.R. (1946). A Statistical Approach to the Problem of the Chronology of Shakespeare's Plays. *Sankhya: Indian J. Statistics*, 7(3), 263-268.

Figure 1.

