

A Preliminary Report on the Set of Symbols Occurring in Engineering Mathematics Texts

Stephen M. Watt
Ontario Research Centre for Computer Algebra
Department of Computer Science
University of Western Ontario
London Ontario, CANADA N6A 5B7
`Stephen.Watt@uwo.ca`

Certain forms of mathematical expression are used more often than others in practice. We propose that a quantitative understanding of actual usage can provide information to improve the accuracy of software for the input of mathematical expressions from scanned documents or handwriting and allow more natural forms of presentation of mathematical expressions by computer algebra systems. Earlier work [1] has examined this question for the diverse set of articles from the mathematics preprint archive `arXiv.org`. That analysis showed the variance between mathematical areas. The present work analyzes a particular mathematical domain more deeply. We have chosen to examine second year university engineering mathematics as taught in North America as the domain. This syllabus typically includes linear algebra, complex analysis, Fourier analysis, vector calculus, and ordinary and partial differential equations. We have analyzed the set of expressions occurring in the most popular textbooks, weighted by popularity. Assuming that early training influences later usage, we take this as a model of the set of mathematical expressions used by the population of North American engineers. We present a preliminary empirical analysis of the individual symbols and of sequences of n symbols (n -grams) occurring in these expressions.

Corpus Selection The first step in our approach was to identify the most popular textbooks in the area of second year engineering mathematics. US college and university bookstore sales for spring for 2006 to fall 2006 show the most demanded texts to be Kreyszig [2] (72%), Greenberg [3] (13%), O’Neil [4] (7%), Jeffrey (5%), Harman (2%). From this we see that three titles account for more than 90% of the textbook use. We therefore built our model based on these three titles.

T_EX Sources For each of the three textbooks, we obtained T_EX sources for all the mathematical expressions, and then constructed MathML from the T_EX. For the texts by Greenberg and O’Neil, the author and publisher (respectively) were highly cooperative and provided the T_EX sources directly. The sources for the text by O’Neil corresponded to the published version in use today. The sources for the text by Greenberg had somewhat diverged from the published text but not so much as to materially affect the analysis in our opinion. For the text by Kreyszig, the publisher and author declined to provide access to the source files. To obtain the mathematical expressions of the text in electronic form, we first scanned the entire book and used the Infty system [5] to produce T_EX. In most cases the T_EX produced had to be edited by hand to correct errors. This was a highly labour intensive activity that spanned several months. In the end we had a T_EX representation for all the mathematical expressions in all three texts.

MathML Conversion Naïve examination of T_EX sources does not give the mathematical expressions of a document. This is for two reasons: The first reason is that typical T_EX document markup makes use of a number of macro packages, as well as author-defined macros. These macros have to be expanded to reveal the mathematical expression. The second reason that T_EX sources do not give *expressions* directly is that the T_EX representation of mathematics is not grouped as required. For example, most authors would write $\$a + b c\$$ rather than $\$a + \{b c\}\$$. We used our T_EX to MathML converter [6, 7] to expand the T_EX macros and properly group the expressions. We then performed our analysis on the resulting MathML. The resulting expressions treated were (for the most part) complete, well formed, and grouped appropriately. We describe the conversion process in more detail elsewhere [1, 8].

Analysis We grouped the chapters of each text into general subject categories (ODEs, PDEs, vector calculus, *etc*) and analyzed the mathematical expressions for each subject/author combination, for each author with subjects combined (weighted by author emphasis), and for each subject with authors combined (weighted by sales volume). In each case, we computed the individual symbol frequencies (normalized to total 1) and n -gram frequencies for $n = 2, 3, 4, 5$. To compute the n -grams, we converted the expressions to strings by traversing the frontier of the expression trees in writing order. The resulting strings were over the alphabet of leaf symbols extended by $\langle \text{sub} \rangle$, $\langle / \text{sub} \rangle$, $\langle \text{sup} \rangle$, $\langle / \text{sup} \rangle$, $\langle \text{frac} / \rangle$ and $\langle \text{root} / \rangle$. These symbols captured transitions from the expression baseline to subscripts and superscripts as well as built up fractions and radicals. The n -grams were then tallied using sliding windows over these strings.

Results Tables 1 and 2 show extracts of the preliminary results of our analysis. Table 1 shows the frequencies of the most commonly occurring symbols in the entire corpus. These are presented with the absolute symbol count for each author and as a percentage of all symbols, weighted by author. The relative weights used were (72, 13, 7). We see that the most popular symbols were common among all the authors, although the rank of the symbols varied somewhat from author to author. The total number of mathematical symbols occurring in the texts were 368,267, 467,044 and 391,602. Table 1 also shows the most commonly occurring symbols for two representative areas. We see that the declining relative frequency is similar between the areas, with a few outlying points (such as z being very popular for complex analysis). This same pattern was observed for all subject areas. The cumulative frequency of symbols is shown in Figure 1 with one curve for each subject and one for the weighted combination. From the log plot it is possible to see that the symbol frequencies follow an approximately exponential distribution.

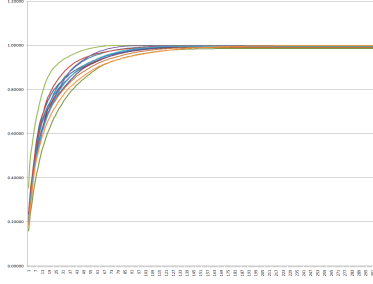
Table 2 shows a preliminary count of the most popular 5-grams for the three corpus authors as well as from two comparison texts. The n -grams have a qualitatively similar declining frequency pattern as the symbols, but this time in a much larger space. The total number of n -grams (for each n) was 479,388 (Kreyszig), 562,297 (Greenberg) and 477,268 (O’Neil). The total number of *different* bigrams was 5,992 (Kreyszig), 7,056 (Greenberg) and 5,442 (O’Neil). The total number of *different* 5-grams was 140,306 (Kreyszig), 146,507 (Greenberg), 126,232 (O’Neil). Figure 1 shows the cumulative frequency for all distinct n -grams occurring in the text by Kreyszig. The highest curve is for $n = 2$ and they are in order to the lowest curve for $n = 5$. We find it remarkable that even though the ranking of the particular n -grams is different for the each author, the cumulative n -gram frequency curves are almost identical from author to author.

By analyzing the population of symbols and n -grams that occur in the corpus, we are able to determine the most popular symbols and n -grams by subject. The exponential drop in number of occurrences, from the highest ranked symbols and n -grams to the lowest, means that a compact database can contain most of the frequently occurring items. Thus applications, even those for portable devices, could use these statistics to guide their recognition.

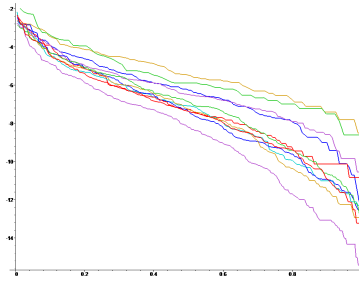
Acknowledgments We thank Michael Greenberg, Peter O’Neil, Prentice-Hall and Thomson-Nelson for the use of their materials. We also thank Robert Lopez and Maplesoft for additional materials. We thank Jeliuzko Polihronov for assistance in gathering the data and Elena Smirnova for work on the n -gram analysis software. This work was supported in part by grants from the NSERC Canada, Microsoft and Maplesoft.

References

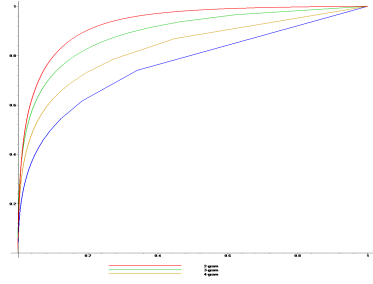
- [1] C.M. So and S.M. Watt, *Determining Empirical Properties of Mathematical Expression Use*, pp. 361-375, Proc. Fourth Int’l Conf. on Mathematical Knowledge Management (MKM 2005), Springer Verlag LNCS 3863.
- [2] Erwin Kreyszig, *Advanced Engineering Mathematics*, 8th edition, John Wiley & Sons 1999.
- [3] Michael Greenberg, *Advanced Engineering Mathematics*, 2nd edition, Prentice Hall 1998.
- [4] Peter O’Neil, *Advanced Engineering Mathematics*, 5th edition, Thomson-Nelson 2003.
- [5] M.Suzuki, F.Tamari, R.Fukuda, S.Uchida, T.Kanahori, *Infty—an Integrated OCR System for Mathematical Documents*, Proceedings of ACM Symposium on Document Engineering 2003, Grenoble, 2003, pp.95-104.
- [6] ORCCA. *On-line TeX to MathML Translator*. <http://www.orcca.on.ca/MathML/texmml/textomml.html>
- [7] S.M.Watt. *Implicit Mathematical Semantics in Conversion between TeX and MathML*. TUGBoat, **23**(1) 2002.
- [8] E. Smirnova and S.M. Watt, *Context-Sensitive Mathematical Character Recognition*. International Conference on Frontiers in Handwriting Recognition (ICFHR 2008), (accepted).



Symbols
X: Symbol rank number
Y: Cumulative frequency



Symbols
X: Symbol rank number
Y: Log frequency



2-, 3-, 4- and 5-grams
X: n -gram rank number
Y: Cumulative frequency

Figure 1: Symbol and n -gram frequencies

All areas combined				
Symbol	Weighted Freq. (%)	Symbol Counts		
		Kreyszig	Greenberg	O'Neil
1	6.16415	24519	23209	20345
2	6.15918	24436	22613	21886
=	5.89883	22906	26202	19275
0	5.13055	20436	19623	16164
(5.08432	18162	26262	27777
)	5.08387	18158	26257	27804
x	4.97402	18271	28243	17918
-	3.82436	14609	15625	17152
+	3.12976	11906	14648	11711
y	2.94812	11400	13191	9996
,	2.53506	9796	12571	6784
,	2.11526	8016	9681	8577
z	1.88590	7447	7238	6593
3	1.87252	7225	7603	7706
□	1.73059	6386	7715	9163
□	1.71003	5771	9800	11446
t	1.62134	6234	4510	10083
.	1.42027	5694	4119	6097
4	1.30925	4926	6522	4874
f				

Complex Analysis	
Symbol	Weighted Freq. (%)
z	11.28007
=	6.19577
)	5.76133
(5.75744
1	5.59297
2	5.21226
-	4.02399
0	3.88584
+	3.71409
i	2.95919
n	2.94910
	2.78406
x	2.45995
f	1.98821
,	1.69837
y	1.60176
π	1.30730
C	1.18192
3	1.13346

PDEs	
Symbol	Weighted Freq. (%)
=	7.22187
x	7.04832
(6.44756
)	6.43967
2	5.54914
0	4.82981
u	4.28608
1	3.35806
n	3.33931
t	3.10607
y	2.63211
-	2.38819
+	2.02753
,	2.00038
c	1.68841
r	1.67920
π	1.66333
f	1.38602
∂	1.32067

Table 1: Most popular symbols, by weighted frequency, for entire corpus and two sample areas

Kreyszig Freq% Sequence	Greenberg Freq% Sequence	O'Neil Freq% Sequence	Lopez Freq% Sequence	MSKit Freq% Sequence
0.00104 (x, y)	0.00142 $\int^{(sub)} 0^{(sup)}$	0.00152 (x, y)	0.00284 00000	0.00442 $\lim^{(sub)} x$
0.00095 $y^{(sup)} \pi^{(sup)}$	0.00130 (x, y)	0.00149 $\int^{(sub)} 0^{(sup)}$	0.00136 (x, y)	0.00406 $im^{(sub)} x \rightarrow$
0.00082 $x^{(sub)} 1^{(sup)} +$	0.00077 $0^{(sup)} \infty^{(sup)}$	0.00106 $y^{(sup)} \pi^{(sup)}$	0.00120 $x^{(sup)} 2^{(sup)} +$	0.00320 $x^{(sup)} 2^{(sup)} +$
0.00081 $f(x) =$	0.00077 $x^{(sup)} 2^{(sup)} +$	0.00102 $0^{(sup)} \infty^{(sup)}$	0.00104 $\int^{(sub)} 0^{(sup)}$	0.00285 $dy^{(frac/)} dx$
0.00080 $\int^{(sub)} 0^{(sup)}$	0.00071 (x, t)	0.00100 $\sum^{(sub)} n = 1$	0.00096 $f(x) =$	0.00200 $f(x) =$
0.00073 $0^{(sup)} \infty^{(sup)}$	0.00070 $1^{(sup)}, \dots$	0.00100 $^{(sub)} n = 1^{(sup)}$	0.00074 $x^{(sup)} 2^{(sup)} -$	0.00200 $\sin(x)$
0.00072 $^{(sup)} \pi^{(sup)} =$	0.00067 $^{(sub)} 1^{(sup)}, \dots$	0.00100 $n = 1^{(sup)}$	0.00070 $(x, y,$	0.00190 $x^{(sup)} 2^{(sup)} -$
0.00072 $x^{(sup)} 2^{(sup)} +$	0.00062 $y(x) =$	0.00094 $1^{(sup)} \infty^{(sup)}$	0.00066 $^{(sup)} 2^{(sup)} + 1$	0.00188 $\ln(x)$
0.00071 $^{(sup)} \pi^{(sup)} +$	0.00060 $^{(sub)} 0^{(sup)} \infty^{(sup)}$	0.00093 $= 1^{(sup)} \infty^{(sup)}$	0.00065 $, \dots,$	0.00154 $2x^{(sup)} 2^{(sup)}$
0.00064 $-z^{(sub)} 0^{(sup)}$	0.00057 $^{(sup)} (x) =$	0.00090 $\sin($	0.00064 $f(x, y$	0.00141 $\cos(x)$
0.00060 \dots	0.00056 $(0) = 0$	0.00086 $x^{(sup)} 2^{(sup)} +$	0.00064 $+y^{(sup)} 2^{(sup)}$	0.00133 $\cos(x)$
0.00057 $^{(sub)} 1^{(sup)} \pi^{(sup)}$	0.00056 $f(x, y$	0.00084 $^{(sup)} \pi^{(sup)} +$	0.00061 x, y, z	0.00131 $x^{(sup)} 3^{(sup)} +$
0.00055 $z^{(sub)} 0^{(sup)}$	0.00052 $^{(sub)} 1^{(sup)} (x$	0.00084 $(-1)^{(sup)}$	0.00060 $, y, z)$	0.00124 $^{(sup)} 2^{(sup)} + 1$
0.00054 $y^{(sub)} 1^{(sup)}$	0.00052 $1^{(sup)} (x)$	0.00082 $\sin(n$	0.00059 $2x^{(sup)} 2^{(sup)}$	0.00122 $\log^{(sub)} a$
0.00054 $y(0) =$	0.00052 $1^{(sup)} \infty^{(sup)}$	0.00076 $y^{(sup)} \pi^{(sup)} =$	0.00058 $^{(sup)} 2^{(sup)} + y$	0.00116 $(x^{(sup)} 2^{(sup)})^{(sub)}$
0.00054 $1^{(sup)} \pi^{(sup)}$	0.00051 x, y, z	0.00073 $y^{(sup)} \pi^{(sup)} +$	0.00058 $\sin($	0.00117 $y^{(frac/)} dx =$
0.00051 $^{(sub)} 2^{(sup)} \pi^{(sup)}$	0.00050 $f(x) =$	0.00072 $, \dots,$	0.00058 $^{(sup)} \cos($	0.00116 $^{(root/)} 2x^{(sup)} 2$
0.00051 $z - z^{(sub)} 0$	0.00049 $(x, y,$	0.00070 (x, t)	0.00057 $^{(sup)} \sin($	0.00116 $(x^{(sup)} 2^{(sup)})^{(sub)}$
0.00050 $, y^{(sub)} 2^{(sup)}$	0.00044 $^{(sub)} n^{(sup)} (x$	0.00068 $-1)^{(sup)} n$	0.00055 $0^{(sup)} \infty^{(sup)}$	0.00116 $du^{(frac/)} dx$

Table 2: Most popular 5-grams