

Survey on Big Data Cleansing

Mika Huttunen
Faculty of science, Helsinki university

ABSTRACT

This paper provides a sample of a \LaTeX document which conforms, somewhat loosely, to the formatting guidelines for ACM SIG Proceedings. It is an *alternate* style which produces a *tighter-looking* paper and was designed in response to concerns expressed, by authors, over page-budgets. It complements the document *Author's (Alternate) Guide to Preparing ACM SIG Proceedings Using $\LaTeX_2\epsilon$ and BibTeX*. This source file has been written with the intention of being compiled under $\LaTeX_2\epsilon$ and BibTeX.

Keywords

big data, data cleansing, data analysis, data transformation

1. INTRODUCTION

The *proceedings* are the records of a conference. ACM seeks to give these conference by-products a uniform, high-quality appearance. To do this, ACM has some rigid requirements for the format of the proceedings documents: there is a specified format (balanced double columns), a specified set of fonts (Arial or Helvetica and Times Roman) in certain specified sizes (for instance, 9 point for body copy), a specified live area (18×23.5 cm [$7'' \times 9.25''$]) centered on the page, specified size of margins (1.9 cm [$0.75''$] top, (2.54 cm [$1''$]) bottom and (1.9 cm [$.75''$]) left and right; specified column width (8.45 cm [$3.33''$]) and gutter size (.83 cm [$.33''$]).

The good news is, with only a handful of manual settings¹, the \LaTeX document class file handles all of this for you.

The remainder of this document is concerned with showing, in the context of an “actual” document, the \LaTeX commands specifically available for denoting the structure of a

¹Two of these, the `\numberofauthors` and `\alignauthor` commands, you have already used; another, `\balancecolumns`, will be used in your very last run of \LaTeX to ensure balanced column heights on the last page.

proceedings paper, rather than with giving rigorous descriptions or explanations of such commands.

2. RELATED WORK

Typically, the body of a paper is organized into a hierarchical structure, with numbered or unnumbered headings for sections, subsections, sub-subsections, and even smaller sections. The command `\section` that precedes this paragraph is part of such a hierarchy.² \LaTeX handles the numbering and placement of these headings for you, when you use the appropriate heading commands around the titles of the headings. If you want a sub-subsection or smaller part to be unnumbered in your output, simply append an asterisk to the command name. Examples of both numbered and unnumbered headings will appear throughout the balance of this sample document.

Because the entire article is contained in the **document** environment, you can indicate the start of a new paragraph with a blank line in your input file; that is why this sentence forms a separate paragraph.

3. GENERAL DATA CLEANSING PROCESS

This section discusses background related to data cleansing. It introduces a general process for data cleansing containing data analysis and data transformation steps.

3.1 Type Changes and *Special Characters*

We have already seen several typeface changes in this sample. You can indicate italicized words or phrases in your text with the command `\textit`.

4. ISSUES WITH BIG DATA

This section discusses problems that arise when data cleansing is done with big data.

5. BIG DATA CLEANSING FRAMEWORKS

This section discusses tools and frameworks designed particularly for handling big data. There are BigDancing and KATARA, but what else is available?

6. WEAK POINTS ON PAPERS

This section discusses weak points of the papers related to the papers current ongoing situation with big data cleansing.

²This is the second footnote. It starts a series of three footnotes that add nothing informational, but just give an idea of how footnotes work and look. It is a wordy one, just so you see how a longish one plays out.

Table 1: Frequency of Special Characters

Non-English or Math	Frequency	Comments
\emptyset	1 in 1,000	For Swedish names
π	1 in 5	Common in math
$\$$	4 in 5	Used in business
Ψ_1^2	1 in 40,000	Unexplained usage

7. OWN IDEAS

Give new ideas/algorithms/experiments on this research problem. This part is very important, because it shows the potential of the author to be an independent innovative researcher. This section can be as long as possible.

Come up with a new idea / algorithm that could be used for dealing with big data and does something better than the existing tools available...

8. CONCLUSION

Summarize the research problem and the main contributions of previous papers. The main weakness of previous works could be also mentioned here. Some future works can be described as well.

This paragraph will end the body of this sample document. Remember that you might still have Acknowledgments or Appendices; brief samples of these follow. There is still the Bibliography to deal with; and we will make a disclaimer about that here: with the exception of the reference to the L^AT_EX book, the citations in this paper are to articles which have nothing to do with the present subject and are used as examples only.

9. REFERENCES

10. EXAMPLES

end the environment with table*, NOTE not table!

As was the case with tables, you may want a figure that spans two columns. To do this, and still to ensure proper “floating” placement of tables, use the environment **figure*** to enclose the figure and its caption. and don’t forget to end the environment with figure*, not figure!

THEOREM 1. *Let f be continuous on $[a, b]$. If G is an antiderivative for f on $[a, b]$, then*

$$\int_a^b f(t)dt = G(b) - G(a).$$

Definition 1. If z is irrational, then by e^z we mean the unique number which has logarithm z :

$$\log e^z = z$$

PROOF. Suppose on the contrary there exists a real number L such that

$$\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = L.$$

Then

$$l = \lim_{x \rightarrow c} f(x) = \lim_{x \rightarrow c} \left[g(x) \cdot \frac{f(x)}{g(x)} \right] = \lim_{x \rightarrow c} g(x) \cdot \lim_{x \rightarrow c} \frac{f(x)}{g(x)} = 0 \cdot L = 0,$$

which contradicts our assumption that $l \neq 0$. \square

Table 2: Some Typical Commands

Command	A Number	Comments
<code>\alignauthor</code>	100	Author alignment
<code>\numberofauthors</code>	200	Author enumeration
<code>\table</code>	300	For tables
<code>\table*</code>	400	For wider tables

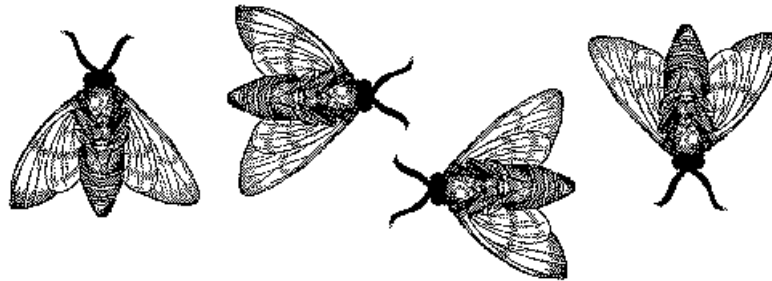


Figure 1: A sample black and white graphic that needs to span two columns of text.

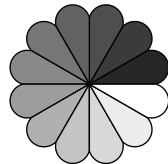


Figure 2: A sample black and white graphic that has been resized with the `includegraphics` command.