

---

# Movie Recommendation based on Collaborative Topic Modeling

---

**Abhishek Bhowmick**  
Department of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213  
abhowmil@andrew.cmu.edu

**Udbhav Prasad**  
Department of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213  
udbhavp@andrew.cmu.edu

**Satwik Kottur**  
Department of Electrical Engineering  
Carnegie Mellon University  
Pittsburgh, PA 15213  
skottur@andrew.cmu.edu

## Abstract

Traditional collaborative filtering relies on reviews provided by viewers in the movie watching community to make recommendations to the user. In this work, we attempt to combine this approach with probabilistic topic modeling techniques to make recommendations that consist not only of movies that are popular in the community, but also those that are similar in content to movies that the user has enjoyed in the past.

## 1 Introduction

Automatic movie/TV show recommendation is an important technology for streaming video services like Netflix, Hulu, HBO, Amazon etc. Not only video streaming services, but also audio/music streaming sites like Spotify, iTunes make heavy use of recommender systems. Indeed, any content management system that has large quantities of information such as usage patterns, metadata, latent structures or the ability to extract such information, can and should make use of recommendation methods to provide services that help find items of interest. In the remaining parts of the paper, we limit ourselves to the domain of movie recommendation, however most of the discussion/analysis can be generally applied to other domains.

### 1.1 Collaborative Filtering and its shortcoming

Traditional collaborative filtering techniques make use of usage patterns, or more specifically, movie reviews. Movie ratings provided by a user are used alongwith similar ratings by other users to build a model that captures the user's preferences. This model is then used to predict movies that the user may have an interest in. However, this method only works when sufficient usage data is available. New content that is available may not be possible to recommend in absence of sufficient usage data. This is known as the 'cold start' problem.

One strategy might be to randomly recommend newly arrived movies to users and record their responses, thus building up usage data. However, such an approach has a few shortcomings. Since an average user likes only a few types of movies, it is more likely than not for the user to give a negative review to a randomly suggested movie. Getting a sufficient number of positive reviews may take a long time through this approach, and also user satisfaction decreases due to the bad

recommendations made by the system (assuming a random recommendation is more likely to be bad than good).

## **1.2 Content Based Recommendation**

The other main approach is content-based recommendation, that addresses the 'cold start' problem of collaborative filtering. A simple way is to make recommendations based on movie metadata such as genre, actors, language etc. However, this approach severely restricts the pool of movies from which new recommendations can be made. It also leads to very predictable results, since recommendations made on the basis of metadata alone resemble the results a user would have got through simple keyword searches.

A much more interesting approach is to identify similarities among movies through latent themes extracted from information such as plot summaries. Topic modeling can be used to describe movies in terms of such latent themes. Such an approach can allow the recommendation system to generalize to new movies that have very little usage data.

## **1.3 Hybrid approach**

## **2 General formatting instructions**

The text must be confined within a rectangle 5.5 inches (33 picas) wide and 9 inches (54 picas) long. The left margin is 1.5 inch (9 picas). Use 10 point type with a vertical spacing of 11 points. Times New Roman is the preferred typeface throughout. Paragraphs are separated by 1/2 line space, with no indentation.

Paper title is 17 point, initial caps/lower case, bold, centered between 2 horizontal rules. Top rule is 4 points thick and bottom rule is 1 point thick. Allow 1/4 inch space above and below title to rules. All pages should start at 1 inch (6 picas) from the top of the page.

For the final version, authors' names are set in boldface, and each name is centered above the corresponding address. The lead author's name is to be listed first (left-most), and the co-authors' names (if different address) are set to follow. If there is only one co-author, list both author and co-author side by side.

Please pay special attention to the instructions in section 4 regarding figures, tables, acknowledgments, and references.

## **3 Headings: first level**

First level headings are lower case (except for first word and proper nouns), flush left, bold and in point size 12. One line space before the first level heading and 1/2 line space after the first level heading.

### **3.1 Headings: second level**

Second level headings are lower case (except for first word and proper nouns), flush left, bold and in point size 10. One line space before the second level heading and 1/2 line space after the second level heading.

#### **3.1.1 Headings: third level**

Third level headings are lower case (except for first word and proper nouns), flush left, bold and in point size 10. One line space before the third level heading and 1/2 line space after the third level heading.

## **4 Citations, figures, tables, references**

These instructions apply to everyone, regardless of the formatter being used.

#### 4.1 Citations within the text

Citations within the text should be numbered consecutively. The corresponding number is to appear enclosed in square brackets, such as [1] or [2]-[5]. The corresponding references are to be listed in the same order at the end of the paper, in the **References** section. (Note: the standard `BIBTEX` style `unsrt` produces this.) As to the format of the references themselves, any style is acceptable as long as it is used consistently.

As submission is double blind, refer to your own published work in the third person. That is, use “In the previous work of Jones et al. [4]”, not “In our previous work [4]”. If you cite your other papers that are not widely available (e.g. a journal paper under review), use anonymous author names in the citation, e.g. an author of the form “A. Anonymous”.

#### 4.2 Footnotes

Indicate footnotes with a number<sup>1</sup> in the text. Place the footnotes at the bottom of the page on which they appear. Precede the footnote with a horizontal rule of 2 inches (12 picas).<sup>2</sup>

#### 4.3 Figures

All artwork must be neat, clean, and legible. Lines should be dark enough for purposes of reproduction; art work should not be hand-drawn. The figure number and caption always appear after the figure. Place one line space before the figure caption, and one line space after the figure. The figure caption is lower case (except for first word and proper nouns); figures are numbered consecutively.

Make sure the figure caption does not get separated from the figure. Leave sufficient space to avoid splitting the figure and figure caption.

You may use color figures. However, it is best for the figure captions and the paper body to make sense if the paper is printed either in black/white or in color.

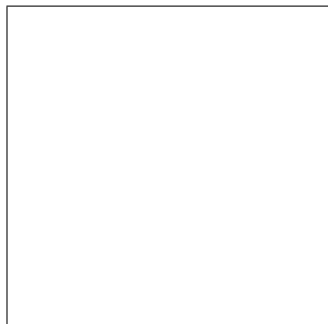


Figure 1: Sample figure caption.

#### 4.4 Tables

All tables must be centered, neat, clean and legible. Do not use hand-drawn tables. The table number and title always appear before the table. See Table 1.

Place one line space before the table title, one line space after the table title, and one line space after the table. The table title must be lower case (except for first word and proper nouns); tables are numbered consecutively.

---

<sup>1</sup>Sample of the first footnote

<sup>2</sup>Sample of the second footnote

Table 1: Sample table title

PART	DESCRIPTION
Dendrite	Input terminal
Axon	Output terminal
Soma	Cell body (contains cell nucleus)

## 5 Final instructions

Do not change any aspects of the formatting parameters in the style files. In particular, do not modify the width or length of the rectangle the text should fit into, and do not change font sizes (except perhaps in the **References** section; see below). Please note that pages should be numbered.

## 6 Preparing PostScript or PDF files

Please prepare PostScript or PDF files with paper size “US Letter”, and not, for example, “A4”. The `-t letter` option on `dvips` will produce US Letter files.

Fonts were the main cause of problems in the past years. Your PDF file must only contain Type 1 or Embedded TrueType fonts. Here are a few instructions to achieve this.

- You can check which fonts a PDF files uses. In Acrobat Reader, select the menu Files>Document Properties>Fonts and select Show All Fonts. You can also use the program `pdf fonts` which comes with `xpdf` and is available out-of-the-box on most Linux machines.
- The IEEE has recommendations for generating PDF files whose fonts are also acceptable for NIPS. Please see <http://www.emfield.org/icuwb2010/downloads/IEEE-PDF-SpecV32.pdf>
- LaTeX users:
  - Consider directly generating PDF files using `pdflatex` (especially if you are a MiKTeX user). PDF figures must be substituted for EPS figures, however.
  - Otherwise, please generate your PostScript and PDF files with the following commands:
 

```
dvips mypaper.dvi -t letter -Ppdf -G0 -o mypaper.ps
ps2pdf mypaper.ps mypaper.pdf
```

 Check that the PDF files only contains Type 1 fonts.
  - `xfig` “patterned” shapes are implemented with bitmap fonts. Use “solid” shapes instead.
  - The `\bbold` package almost always uses bitmap fonts. You can try the equivalent AMS Fonts with command
 

```
\usepackage[psamsfonts]{amssymb}
```

 or use the following workaround for reals, natural and complex:
 

```
\newcommand{\RR}{\mathbb{R}} %real numbers
\newcommand{\Nat}{\mathbb{N}} %natural numbers
\newcommand{\CC}{\mathbb{C}} %complex numbers
```
  - Sometimes the problematic fonts are used in figures included in LaTeX files. The ghostscript program `eps2eps` is the simplest way to clean such figures. For black and white figures, slightly better results can be achieved with program `potrace`.
- MSWord and Windows users (via PDF file):
  - Install the Microsoft Save as PDF Office 2007 Add-in from <http://www.microsoft.com/downloads/details.aspx?displaylang=en&familyid=4d951911-3e7e-4ae6-b059-a2e79ed87041>

- Select “Save or Publish to PDF” from the Office or File menu
- MSWord and Mac OS X users (via PDF file):
  - From the print menu, click the PDF drop-down box, and select “Save as PDF...”
- MSWord and Windows users (via PS file):
  - To create a new printer on your computer, install the AdobePS printer driver and the Adobe Distiller PPD file from <http://www.adobe.com/support/downloads/detail.jsp?ftpID=204> *Note:* You must reboot your PC after installing the AdobePS driver for it to take effect.
  - To produce the ps file, select “Print” from the MS app, choose the installed AdobePS printer, click on “Properties”, click on “Advanced.”
  - Set “TrueType Font” to be “Download as Softfont”
  - Open the “PostScript Options” folder
  - Select “PostScript Output Option” to be “Optimize for Portability”
  - Select “TrueType Font Download Option” to be “Outline”
  - Select “Send PostScript Error Handler” to be “No”
  - Click “OK” three times, print your file.
  - Now, use Adobe Acrobat Distiller or ps2pdf to create a PDF file from the PS file. In Acrobat, check the option “Embed all fonts” if applicable.

If your file contains Type 3 fonts or non embedded TrueType fonts, we will ask you to fix it.

## 6.1 Margins in LaTeX

Most of the margin problems come from figures positioned by hand using `\special` or other commands. We suggest using the command `\includegraphics` from the `graphicx` package. Always specify the figure width as a multiple of the line width as in the example below using `.eps` graphics

```
\usepackage[dvips]{graphicx} ...
\includegraphics[width=0.8\linewidth]{myfile.eps}
```

or

```
\usepackage[pdftex]{graphicx} ...
\includegraphics[width=0.8\linewidth]{myfile.pdf}
```

for `.pdf` graphics. See section 4.4 in the `graphics` bundle documentation (<http://www.ctan.org/tex-archive/macros/latex/required/graphics/grfguide.ps>)

A number of width problems arise when LaTeX cannot properly hyphenate a line. Please give LaTeX hyphenation hints using the `\-` command.

## Acknowledgments

Use unnumbered third level headings for the acknowledgments. All acknowledgments go at the end of the paper. Do not include acknowledgments in the anonymized submission, only in the final paper.

## References

References follow the acknowledgments. Use unnumbered third level heading for the references. Any choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font size to ‘small’ (9-point) when listing the references. **Remember that this year you can use a ninth page as long as it contains *only* cited references.**

[1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauero, D. S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609-616. Cambridge, MA: MIT Press.

- [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the GEneral NEural Simulation System*. New York: TELOS/Springer-Verlag.
- [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.