

# Recommendation Systems for Movies with Sentiment Analysis using Neural Networks

Colby Wise  
Columbia University  
cjlw2165@columbia.edu

Michael Alvarino  
Columbia University  
maa2282@columbia.edu

Richard Dewey  
Columbia University  
rld2126@columbia.edu

## Abstract

*In this research paper we apply the methodology outlined in the Google Research paper "Wide and Deep Learning for Recommender Systems" to predict user ratings of movies. Unlike the Google Wide and Deep model based on only categorical features, our approach utilizes user sentiment analysis by incorporating text-based movie reviews to enhance model prediction accuracy.*

## 1. Introduction

Our approach is unique in that unlike the Google Wide and Deep model based on only categorical features, our approach incorporates text-based movie reviews to enhance the our model's prediction accuracy. Our initial milestone recreated the Google model similarly using only categorical features but with our data. This will serve as a baseline for model comparison and evaluation. Going forward, we will add in text based reviews in our model then experiment with convolutional and recurrent neural network architectures to learn the sentiment of reviews with the goal of more accurately predicting movie ratings.

Source Code:: Wide & Deep Learning for Movie Recommendation Systems

### 1.1. Related Work

The literature on recommendation systems is deep and until recently has traditionally focused on well-known matrix factorization algorithms, such as collaborative filtering, as popularized by Netflix and Spotify [Lee, et. al., 2012]. The benefits of collaborative filtering is that it is relatively easy to implement and computationally efficient given most similarity measure i.e. cosine similarity is matrix multiplication. Conversely, the draw backs in real-world applications have significant implications: namely CF suffers from the "cold-start" problem in that it requires user ratings to make predictions. For instance, for new users without prior history it's hard to predict ratings/recommendations

thus models have poor generalization [Balakrishnan, Dixit, 2016]

Recent research has focused on using Deep Learning, specifically convolutional and recurrent neural networks to help solve the generalization problem. Specifically, we will be focusing on expanding the Google Wide and Deep model thus our primary reference will be the original [Cheng, Koc, et. al. 2016] paper. Google provided a very basic tutorial outlining the general TensorFlow implementation of a generic Wide and Deep model. According to our knowledge no specific source code referencing the paper is open-source, however there are a few similar code-bases which we will reference open-source examples. Additional related papers that we referenced for constructing the Deep portion of our network include DeepFM [Guo, Tang, et. al. 2017] which similar to us, combine a matrix factorization with a CNN architecture; and DeepCoNN [Zheng, Noroozi, Yu, 2017] which uses parallel CNN architectures to predict ratings from users text reviews.

### 1.2. Problem Formulation

*In order to make better movie recommendations how can we more accurately predict a users rating for an unseen movie based on what the user has previously seen? Furthermore, can we improve generalization of our model when information on a users past movie ratings is limited?*

These two questions are the crux of our problem. Prior to deep learning, standard approaches for recommendation systems (*RecSys*) used collaborative filtering which relies on decomposing users, items (i.e. movies), and ratings into latent feature matrices. Then the weights of these matrices are used to predict a rating a user would give for an item. One common method includes using the cosine similarity measure between all pairs of movies that users have rated:

Where,

$m_i$  and  $m_j$  refer to movie1/movie2 and denote vectors of ratings from users have rated both movies:

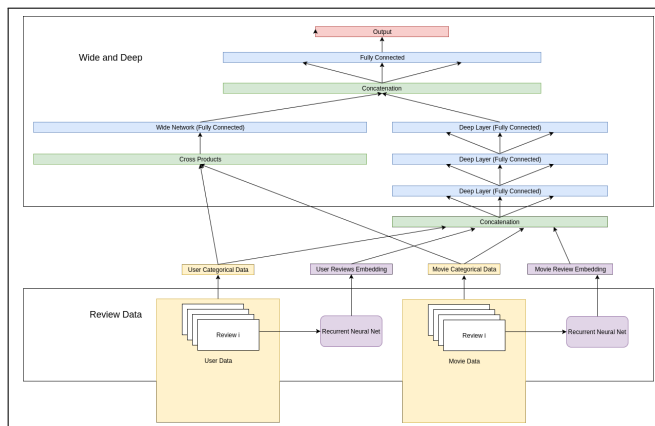


Figure 1. Example of caption. It is set in Roman so that mathematics (always set in Roman:  $B \sin A = A \sin B$ ) may be included without an ugly clash.

$$\text{sim}(m_i, m_j) = \cos(\theta) = \frac{\vec{m}_i \bullet \vec{m}_j}{\|\vec{m}_i\|_2 \times \|\vec{m}_j\|_2}$$

This yields a movie-to-movie similarity matrix of dimensions  $M \times M$  with ones along the diagonal. Thus, the predicted rating for movie  $m_2$  for user1 would be calculated using similarity measures between  $(m_2, m_1)$  and  $(m_2, m_3)$  weighted by the respective ratings for  $m_1$  and  $m_3$ .

From the above formulation we can clearly see the major disadvantage of this approach: sparsity. When ratings are limited the movie-to-movie matrix is mainly zeros thus limiting predictive ability. Google’s *Wide & Deep Model* attempts to work around this by training two networks in parallel. Both are discussed in more detail below, but the **Wide network** is a generalized linear model of input features and transformed features. The **Deep network** is a three layer feed-forward neural network with non-linear *ReLU* activation layers.

### 1.3. Methodology

#### Design

#### Architecture

Please number all of your sections and displayed equations. It is important for readers to be able to refer to any particular equation. Just because you didn’t refer to it in the text doesn’t mean some future reader might not need to refer to it. It is cumbersome to have to use circumlocutions like “the equation second from the top of page 3 column 1”. (Note that the ruler will not be present in the final copy, so is not an alternative to equation numbers). All authors will benefit from reading Mermin’s description of how to write mathematics: <http://www.pamitc.org/documents/mermin.pdf>.

### 1.4. Preliminary Results

Many authors misunderstand the concept of anonymizing for blind review. Blind review does not mean that one must remove citations to one’s own work—in fact it is often impossible to review a paper unless the previous citations are known and available.

Blind review means that you do not use the words “my” or “our” when citing previous work. That is all. (But see below for techreports.)

Saying “this builds on the work of Lucy Smith [1]” does not say that you are Lucy Smith; it says that you are building on her work. If you are Smith and Jones, do not say “as we show in [7]”, say “as Smith and Jones show in [7]” and at the end of the paper, include reference 7 as you would any other cited work.

An example of a bad paper just asking to be rejected:

An analysis of the frobnicatable foo filter.

In this paper we present a performance analysis of our previous paper [1], and show it to be inferior to all previously known methods. Why the previous paper was accepted without this analysis is beyond me.

[1] Removed for blind review

An example of an acceptable paper:

An analysis of the frobnicatable foo filter.

In this paper we present a performance analysis of the paper of Smith *et al.* [1], and show it to be inferior to all previously known methods. Why the previous paper was accepted without this analysis is beyond me.

[1] Smith, L and Jones, C. “The frobnicatable foo filter, a fundamental contribution to human knowledge”. *Nature* 381(12), 1-213.

If you are making a submission to another conference at the same time, which covers similar or overlapping material, you may need to refer to that submission in order to explain the differences, just as you would if you had previously published related work. In such cases, include the anonymized parallel submission [?] as additional material and cite it as

[1] Authors. “The frobnicatable foo filter”, F&G 2014 Submission ID 324, Supplied as additional material [fg324.pdf](#).

Finally, you may feel you need to tell the reader that more details can be found elsewhere, and refer them to a technical report. For conference submissions, the paper must stand on its own, and not *require* the reviewer to go

to a techreport for further details. Thus, you may say in the body of the paper “further details may be found in [?]”. Then submit the techreport as additional material. Again, you may not assume the reviewers will read this material.

Sometimes your paper is about a problem which you tested using a tool which is widely known to be restricted to a single institution. For example, let’s say it’s 1969, you have solved a key problem on the Apollo lander, and you believe that the CVPR70 audience would like to hear about your solution. The work is a development of your celebrated 1968 paper entitled “Zero-g frobnication: How being the only people in the world with access to the Apollo lander source code makes us a wow at parties”, by Zeus *et al.*

You can handle this paper like any other. Don’t write “We show how to improve our previous work [Anonymous, 1968]. This time we tested the algorithm on a lunar lander [name of lander removed for blind review]”. That would be silly, and would immediately identify the authors. Instead write the following:

We describe a system for zero-g frobnication. This system is new because it handles the following cases: A, B. Previous systems [Zeus *et al.* 1968] didn’t handle case B properly. Ours handles it by including a foo term in the bar integral.

...

The proposed system was integrated with the Apollo lunar lander, and went all the way to the moon, don’t you know. It displayed the following behaviours which show how well we solved cases A and B: ...

As you can see, the above text follows standard scientific convention, reads better than the first version, and does not explicitly name you as the authors. A reviewer might think it likely that the new paper was written by Zeus *et al.*, but cannot make any decision based on that guess. He or she would have to be sure that no other authors could have been contracted to solve problem B.

## FAQ

**Q:** Are acknowledgements OK?

**A:** No. Leave them for the final copy.

**Q:** How do I cite my results reported in open challenges?

**A:** To conform with the double blind review policy, you can report results of other challenge participants together with your results in your paper. For your results, however, you should not identify yourself and should not mention your participation in the challenge. Instead present your results referring to the method proposed in your paper and draw conclusions based on the experimental comparison to other results.

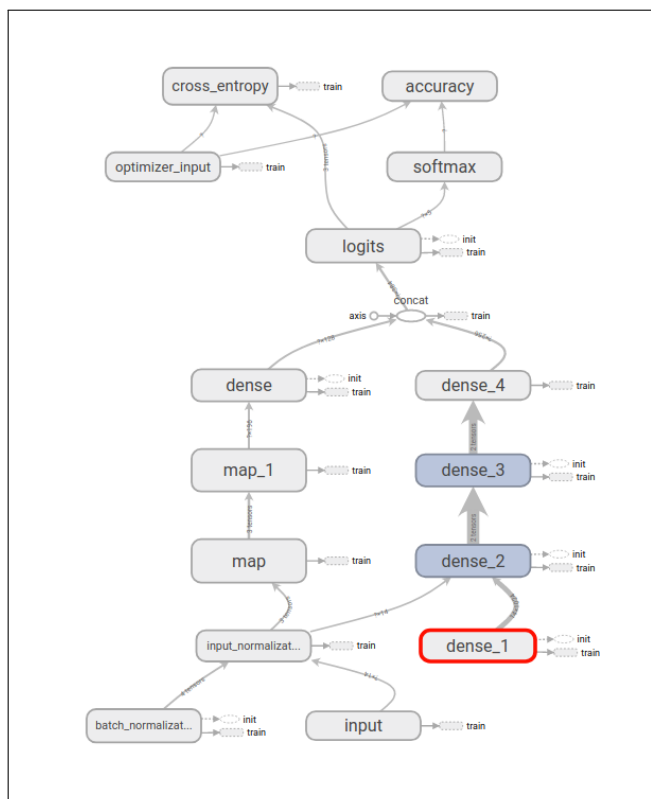


Figure 2. Example of caption. It is set in Roman so that mathematics (always set in Roman:  $B \sin A = A \sin B$ ) may be included without an ugly clash.

## 1.5. Miscellaneous

Compare the following:

$\$conf\_a\$$   $conf_a$   
 $\$\mathit{conf}\_a\$$   $conf_a$

See The T<sub>E</sub>Xbook, p165.

The space after *e.g.*, meaning “for example”, should not be a sentence-ending space. So *e.g.* is correct, *e.g.* is not. The provided `\eg` macro takes care of this.

When citing a multi-author paper, you may save space by using “et alia”, shortened to “*et al.*” (not “*et. al.*” as “*et*” is a complete word.) However, use it only when there are three or more authors. Thus, the following is correct: “Frobnication has been trendy lately. It was introduced by Alpher [?], and subsequently developed by Alpher and Fotheringham-Smythe [?], and Alpher *et al.* [?].”

This is incorrect: “... subsequently developed by Alpher *et al.* [?] ...” because reference [?] has just two authors. If you use the `\etal` macro provided, then you need not worry about double periods when used at the end of a sentence as in Alpher *et al.*

For this citation style, keep multiple citations in numerical (not chronological) order, so prefer [?, ?, ?] to [?, ?, ?].



Figure 3. Example of a short caption, which should be centered.

## 2. Formatting your paper

All text must be in a two-column format. The total allowable width of the text area is  $6\frac{7}{8}$  inches (17.5 cm) wide by  $8\frac{7}{8}$  inches (22.54 cm) high. Columns are to be  $3\frac{1}{4}$  inches (8.25 cm) wide, with a  $\frac{5}{16}$  inch (0.8 cm) space between them. The main title (on the first page) should begin 1.0 inch (2.54 cm) from the top edge of the page. The second and following pages should begin 1.0 inch (2.54 cm) from the top edge. On all pages, the bottom margin should be 1-1/8 inches (2.86 cm) from the bottom edge of the page for 8.5 × 11-inch paper; for A4 paper, approximately 1-5/8 inches (4.13 cm) from the bottom edge of the page.

### 2.1. Margins and page numbering

All printed material, including text, illustrations, and charts, must be kept within a print area 6-7/8 inches (17.5 cm) wide by 8-7/8 inches (22.54 cm) high. Page numbers should be in footer with page numbers, centered and .75 inches from the bottom of the page and make it start at the correct page number rather than the 4321 in the example. To do this fine the line (around line 23)

```
%\ifcvprfinal\pagestyle{empty}\fi
\setcounter{page}{4321}
```

where the number 4321 is your assigned starting page.

Make sure the first page is numbered by commenting out the first page being empty on line 46

```
%\thispagestyle{empty}
```

### 2.2. Type-style and fonts

Wherever Times is specified, Times Roman may also be used. If neither is available on your word processor, please use the font closest in appearance to Times to which you have access.

**MAIN TITLE.** Center the title 1-3/8 inches (3.49 cm) from the top edge of the first page. The title should be

in Times 14-point, boldface type. Capitalize the first letter of nouns, pronouns, verbs, adjectives, and adverbs; do not capitalize articles, coordinate conjunctions, or prepositions (unless the title begins with such a word). Leave two blank lines after the title.

**AUTHOR NAME(s)** and **AFFILIATION(s)** are to be centered beneath the title and printed in Times 12-point, non-boldface type. This information is to be followed by two blank lines.

The **ABSTRACT** and **MAIN TEXT** are to be in a two-column format.

**MAIN TEXT.** Type main text in 10-point Times, single-spaced. Do NOT use double-spacing. All paragraphs should be indented 1 pica (approx. 1/6 inch or 0.422 cm). Make sure your text is fully justified—that is, flush left and flush right. Please do not place any additional blank lines between paragraphs.

Figure and table captions should be 9-point Roman type as in Figures ?? and ??. Short captions should be centred.

Callouts should be 9-point Helvetica, non-boldface type. Initially capitalize only the first word of section titles and first-, second-, and third-order headings.

**FIRST-ORDER HEADINGS.** (For example, **1. Introduction**) should be Times 12-point boldface, initially capitalized, flush left, with one blank line before, and one blank line after.

**SECOND-ORDER HEADINGS.** (For example, **1.1. Database elements**) should be Times 11-point boldface, initially capitalized, flush left, with one blank line before, and one after. If you require a third-order heading (we discourage it), use 10-point Times, boldface, initially capitalized, flush left, preceded by one blank line, followed by a period and your text on the same line.

Method	Frobnability
Theirs	Frumpy
Yours	Frobbly
Ours	Makes one's heart Frob

Table 1. Results. Ours is better.

Please direct any questions to the production editor in charge of these proceedings at the IEEE Computer Society Press: Phone (714) 821-8380, or Fax (714) 761-1784.

### 2.3. Footnotes

Please use footnotes<sup>1</sup> sparingly. Indeed, try to avoid footnotes altogether and include necessary peripheral observations in the text (within parentheses, if you prefer, as in this sentence). If you wish to use a footnote, place it at the bottom of the column on the page on which it is referenced. Use Times 8-point type, single-spaced.

### 2.4. References

List and number all bibliographical references in 9-point Times, single-spaced, at the end of your paper. When referenced in the text, enclose the citation number in square brackets, for example [?]. Where appropriate, include the name(s) of editors of referenced books.

### 2.5. Illustrations, graphs, and photographs

All graphics should be centered. Please ensure that any point you wish to make is resolvable in a printed copy of the paper. Resize fonts in figures to match the font in the body text, and choose line widths which render effectively in print. Many readers (and reviewers), even of an electronic copy, will choose to print your paper in order to read it. You cannot insist that they do otherwise, and therefore must not assume that they can zoom in to see tiny details on a graphic.

When placing figures in  $\text{\LaTeX}$ , it's almost always best to use `\includegraphics`, and to specify the figure width as a multiple of the line width as in the example below

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

### 2.6. Color

Please refer to the author guidelines on the CVPR 2018 web page for a discussion of the use of color in your document.

## 3. Final copy

You must include your signed IEEE copyright release form when you submit your finished paper. We MUST have this form before your paper can be published in the proceedings.

---

<sup>1</sup>This is what a footnote looks like. It often distracts the reader from the main flow of the argument.