

UNIVERSIDADE FEDERAL DO RIO GRANDE DO SUL  
INSTITUTO DE INFORMÁTICA  
CURSO DE CIÊNCIA DA COMPUTAÇÃO

ADRIANO CARNIEL BENIN

**A comparison of recommender systems for  
crowdfunding projects**

Work presented in partial fulfillment  
of the requirements for the degree of  
Bachelor in Computer Science

Advisor: Prof. Dr. Bruno Castro da Silva

Porto Alegre  
March 2018

UNIVERSIDADE FEDERAL DO RIO GRANDE DO SUL

Reitor: Prof. Rui Vicente Oppermann

Vice-Reitora: Prof<sup>a</sup>. Jane Fraga Tutikian

Pró-Reitor de Graduação: Prof. Wladimir Pinheiro do Nascimento

Diretora do Instituto de Informática: Prof<sup>a</sup>. Carla Maria Dal Sasso Freitas

Coordenador do Curso de Ciência de Computação: Prof. Raul Fernando Weber

Bibliotecária-chefe do Instituto de Informática: Beatriz Regina Bastos Haro

## ABSTRACT

Este documento é um exemplo de como formatar documentos para o Instituto de Informática da UFRGS usando as classes L<sup>A</sup>T<sub>E</sub>X disponibilizadas pelo UTUG. Ao mesmo tempo, pode servir de consulta para comandos mais genéricos. *O texto do resumo não deve conter mais do que 500 palavras.*

**Keywords:** UFRGS. recommender systems. AI. crowdfunding.

## Using L<sup>A</sup>T<sub>E</sub>X to Prepare Documents at II/UFRGS

### RESUMO

This document is an example on how to prepare documents at II/UFRGS using the L<sup>A</sup>T<sub>E</sub>X classes provided by the UTUG. At the same time, it may serve as a guide for general-purpose commands. *The text in the abstract should not contain more than 500 words.*

**Palavras-chave:** Electronic document preparation. L<sup>A</sup>T<sub>E</sub>X. ABNT. UFRGS.

## LIST OF FIGURES

Figure 1.1 Exemplo de figura importada de um arquivo e também exemplo de caption muito grande que ocupa mais de uma linha na Lista de Figuras .....	9
Figure 5.1 GbTree example.....	15

## **LIST OF TABLES**

## **LIST OF ABBREVIATIONS AND ACRONYMS**

SMP	Symmetric Multi-Processor
CBF	Content-based filtering
CF	Collaborative filtering
SPMD	Single Program Multiple Data
ABNT	Associação Brasileira de Normas Técnicas

## CONTENTS

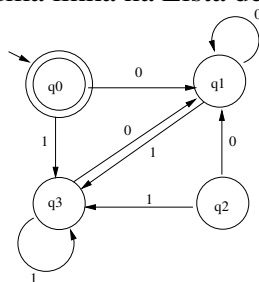
<b>1 INTRODUÇÃO .....</b>	<b>9</b>
<b>2 RECOMMENDATION ALGORITHMS .....</b>	<b>10</b>
2.1 Collaborative filtering.....	10
2.2 Content-based filtering .....	11
<b>3 ABOUT CATARSE .....</b>	<b>13</b>
<b>4 RECOMMENDING PROJECTS.....</b>	<b>14</b>
4.1 Data preparation.....	14
<b>5 GRADIENT BOOSTING TREES .....</b>	<b>15</b>
5.1 Objective function.....	15
<b>REFERENCES.....</b>	<b>17</b>



## 1 INTRODUÇÃO

No início dos tempos, Donald E. Knuth criou o  $\text{T}_{\text{E}}\text{X}$ . Algum tempo depois, Leslie Lamport criou o  $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$ . Graças a eles, não somos obrigados a usar o Word nem o Libre-Office.

Figure 1.1: Exemplo de figura importada de um arquivo e também exemplo de caption muito grande que ocupa mais de uma linha na Lista de Figuras



Fonte: Os Autores

## 2 RECOMMENDATION ALGORITHMS

Recommendation Algorithms are widely used in the industry today to provide useful suggestions to end-users in a completely automated manner. They are ubiquitous in modern e-commerce Web sites (SCHAFER; KONSTAN; RIEDL, 2001), where new products can be recommended based on a customer's interests and preferences, and in many other fields such as movies (Netflix) and music (Spotify). Its importance can't be overstated: the effectiveness of targeted recommendations, as measured by click-through and conversion rates, far exceed those of untargeted content (LINDEN; SMITH; YORK, 2003).

The basic idea behind any recommender system is to obtain a utility function to estimate a user preferences towards an item. The meaning of this function will differ for each context; it could mean how likely a user will want to watch a specific movie or listen to a song, or the likelihood of buying a particular product. In our case, the goal is to find projects the customer is most likely to back given his backing history and other characteristics.

Recommender systems can be broadly divided into two categories: content-based methods and collaborative filtering methods (RAKESH; LEE; REDDY, 2016). The former utilizes the content features of users or items in order to recommend items to users. In collaborative filtering methods, user ratings and other data are used in order to calculate similarities between users or items, which are then ranked to show the most relevant recommendations. These two methods are sometimes combined into what is known as Hybrid Recommender Systems.

### 2.1 Collaborative filtering

Collaborative filtering is based on the principle that similar users will share similar interests. In the traditional approach, each customer is represented by a  $N$ -dimensional vector of items, where  $N$  stands for the number of available items, and each vector component corresponds to the user rating of the item. For large databases, CF can be prohibitively computationally expensive. Its worst case performance is  $O(MN)$  where  $M$  is the number of customers and  $N$  is the number of items (LINDEN; SMITH; YORK, 2003), but this problem can be generally alleviated since in most cases the customer vector is extremely sparse. Item-to-item CF is another method proposed by Amazon that tries

to minimize these scaling issues by focusing on item instead of user similarity. Each item purchased by the customer is compared to other items in the dataset in order to calculate a similarity metric. The algorithm is shown bellow:

```

for each item in product catalog, I1 do
    | for each customer C who purchased I1 do
    | | for item I2 purchased by customer C do
    | | | Record that a customer purchased I1 and I2;
    | | end
    | end
    | for each item I2 do
    | | Compute the similarity between I1 and I2;
    | end
end

```

#### **Algorithm 1:** Item-to-item CF

There are many ways to compute the similarity between items, one common method is the cosine measure defined as follows:

$$\text{similarity}(x, y) = \cos(x, y) = \frac{x \cdot y}{\|x\| \cdot \|y\|}$$

In either case, for this to work large amounts of data about the user is required. This is known as the cold start problem: new users who haven't rated many items yet will have a reduced recommendation quality. On the other hand, no information about the item itself is needed, making CF specially applicable to collections of hard-to-analyze items such as movies.

## **2.2 Content-based filtering**

In this method a description of each item is constructed using structured data or some sort of item presentation algorithm such as Latent Dirichlet Allocation or TF-IDF. This representation is then compared to the user profile and the best-matching items are recommended. The user profile can consist of many different types of information: a history of the user's interaction with the system such as page views, searches and purchases is often used to train the model without any explicit user input. Some systems may require the user to explicitly state his interests, however it may be very hard to get users to make this effort, rendering this approach very limited in practice.

As CBF focus on item rather than user similarity, it avoids the cold start problem

since little information about user preference is needed. However, since only similar items to those already rated by the user will be considered, CBF approaches tend to suffer from over-specialization (IAQUINTA et al., 2008). This is known as the serendipity problem. Another limitation of CBF algorithms is that items are required to contain enough information in order to distinguish items the user likes from items the user doesn't like. For example, a dataset of songs where only the song name is available would not be enough to make good predictions based on this content only, however this would pose no problem for collaborative filtering methods which rely on user similarity.

### **3 ABOUT CATARSE**

Launched in January 2011, Catarse was the first crowdfunding platform for creative projects in Brazil. With over 7000 successfully financed projects raising R\$77 millions from 480.000 people, it's currently the largest national platform of its kind. It works similarly to most crowdfunding platforms: the project owner presents his or her idea and specifies the required investment as well as the cutoff date for the project, while offering rewards for those who back it. Projects are divided into 3 main categories: all-or-nothing, flexible and recurrent. In the first type, projects are available for backing up to 60 days and the project owner only receives the raised amount if the project's goal is met, otherwise all the money is returned to its original backers. On flexible projects the owner receives the raised amount whether the goal is reached or not. Recurrent projects are subscription based and the owner can collect the money monthly. This work will only focus on the first two types of projects.

## 4 RECOMMENDING PROJECTS

Catarse’s dataset is unique in a variety of ways that makes collaborative filtering techniques hard to apply. Firstly, unlike regular e-commerce Web sites such as Amazon where products stay available for many years, crowdfunding projects have a predetermined cutoff date, after which it’s no longer possible to make a pledge. It makes no sense to recommend expired projects, severely limiting our training data to online projects only. Another challenge is the fact that the majority of Catarse’s users only back one project, making it very hard to get enough data for CF methods to work properly. Our only choice is then to use content-based methods, this allows us to train our model with the whole dataset to extract backer-project features, and later use this model to search for online projects with the highest backing probability for the current user.

### 4.1 Data preparation

The dataset  $D$  was obtained directly from Catarse’s production database. Each training example consists of a project-backer pair  $(p, b)$  with 29 dimensions divided as follows:

- *project features(10)*: category, mode, location, recommended by admin, duration, goal, budget length, description length, total contributions, number of projects by the owner
- *user features(19)*: location, number of backed projects in each category

In order to trim irrelevant data, we remove canceled and draft projects as well as projects with no backers from the dataset. The final cardinality of  $D$  was 1103836. In order to balance our dataset, we create 1103836 negative examples by randomly combining users with projects they haven’t backed. test auc k-5 0.990905 std 0.000287 Accuracy : 0.9645 final: train-auc:0.983985

## 5 GRADIENT BOOSTING TREES

In recent years, tree boosting has become an increasingly popular method and has been shown to give state-of-the-art results for many classification problems (LI, ). As in any supervised learning method, our objective is to train a model to predict a target variable  $y_i$  given features  $x_i$ . Boosting methods are characterized by combining many weak learners into a strong one in an iterative fashion. In this case, our model is an ensemble of trees, more specifically a set of classification and regression trees (CART). Unlike decision trees, where the leaves contain decision values, the leaves on gradient boosting trees contain a score  $s \in R$ . Multiple simple trees are then constructed and the prediction of each tree is summed up to get the final score 5.1. Our model can then be defined as

$$\hat{y}_i = \sum_{k=1}^K f_k(x_i), f_k \in \mathcal{F}$$

where  $K$  is the number of trees,  $f$  is a function in the functional space  $\mathcal{F}$ , and  $\mathcal{F}$  is the set of all possible CARTs (CHEN; GUESTRIN, 2016). Our objective function can be written as

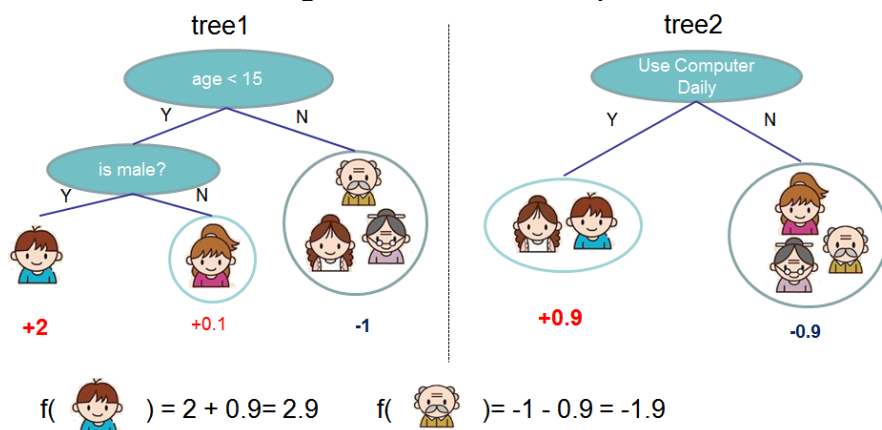
$$obj(\theta) = \sum_i^n l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

### 5.1 Objective function

Our objective function is used to measure the performance of our model for each set of parameters. We define it as the sum of the training loss function  $L$  with the regularization term  $\Omega$ .

$$obj(\theta) = L(\theta) + \Omega(\theta)$$

Figure 5.1: GbTree example



Source: XGBoost official site

In this work we define our loss function as the logistic loss function for logistic regression.

$$L(\theta) = \sum_i [y_i \ln(1 + e^{-\hat{y}_i}) + (1 - y_i) \ln(1 + e^{\hat{y}_i})]$$

And our regularization function is defined as follows:

$$\Omega(f) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

Where  $w$  is the vector of scores on leaves and  $T$  is the number of leaves.

After defining our objective function, we optimize it in order to train our model.



## REFERENCES

- CHEN, T.; GUESTIN, C. XGBoost. In: **Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD 16**. [S.l.]: ACM Press, 2016.
- IAQUINTA, L. et al. Introducing serendipity in a content-based recommender system. In: **2008 Eighth International Conference on Hybrid Intelligent Systems**. [S.l.]: IEEE, 2008.
- LI, P. Robust logitboost and adaptive base class (abc) logitboost.
- LINDEN, G.; SMITH, B.; YORK, J. Amazon.com recommendations: item-to-item collaborative filtering. **IEEE Internet Computing**, Institute of Electrical and Electronics Engineers (IEEE), v. 7, n. 1, p. 76–80, jan 2003.
- RAKESH, V.; LEE, W.-C.; REDDY, C. K. Probabilistic group recommendation model for crowdfunding domains. In: **Proceedings of the Ninth ACM International Conference on Web Search and Data Mining - WSDM 16**. [S.l.]: ACM Press, 2016.
- SCHAFER, J. B.; KONSTAN, J. A.; RIEDL, J. E-commerce recommendation applications. In: **Applications of Data Mining to Electronic Commerce**. [S.l.]: Springer US, 2001. p. 115–153.