

Modeling Dynamic User Interests: A Neural Matrix Factorization Approach

Paramveer S. Dhillon, Sinan Aral

MARKETING SCIENCE 2021

2023 年 10 月 23 日



- 1 Introduction
- 2 Related Work
- 3 Model
- 4 Data
- 5 Results
- 6 Discussion

1 Introduction

2 Related Work

3 Model

4 Data

5 Results

6 Discussion

Background

- ◎ 消费者数据（消费者行为数据和 UGC）被用于营销领域的各类研究中；
- ◎ 三大挑战：
 - 高维稀疏的非结构化数据，使用传统方法进行统计推断往往比较困难；
 - 数据生成过程本身的非平稳性（动态性）；
 - 不同于市场营销研究中常见的购买数据建模研究，用户内容消费数据的建模更加复杂，内容消费的种类总是在不断增加，而且几乎没有重复的“动机”。

Research Question and Objective

- ◎ 本文旨在解决上述三大挑战，提出了一种基于神经网络的矩阵分解模型，用于建模用户的动态兴趣；
- ◎ 模型结果具有可解释性，同时还保留了神经网络的灵活性。

- 1 Introduction
- 2 Related Work**
- 3 Model
- 4 Data
- 5 Results
- 6 Discussion

Related Work

- ◎ 消费者行为建模
 - 使用 Topic Model 对用户消费的文本数据进行建模，从而生成用户画像 (Trusov et al. 2016).
- ◎ 消费者偏好演变及其对各种营销变量敏感性
 - 指数平滑模拟品牌偏好的演变 (Guadagni and Little 1983);
 - 高斯过程模拟消费者偏好的动态变化 (Dew et al. 2020);
 - 使用粒子滤波对广告和销售等其他营销变量之间的非线性关系进行建模 (Bruce 2008)
- ◎ 市场营销领域使用机器学习方法研究客户兴趣 (Netzer et al. 2012, Tirunillai and Tellis 2014, B"uschken and Allenby 2016, Liu and Toubia 2018, Timoshenko and Hauser 2019)
- ◎ 矩阵分解

- 1 Introduction
- 2 Related Work
- 3 Model**
- 4 Data
- 5 Results
- 6 Discussion

Notation

- ⊙ n users, T time slices;
- ⊙ Content consumed by user i in time period t and user i 's unique identity:

$$\mathbf{z}_i^t = \left[\mathbf{x}_i^t; \mathbf{a}_i \right]$$

$p \times 1$ one-hot vector $n \times 1$ vector

- ⊙ All content consumed by all users during the whole time slices:

$$\{\mathbf{Z}\}^{t=1:\tau}$$

$n \times (n + p)$

Matrix Factorization for Modeling Users' Content Interests I

- Approximate the data matrix using the user and content factors

$$z_{ij}^t \approx \underbrace{v_j^\top}_{K \times (n+p) \text{ item vector}} \underbrace{u_i^t}_{K \times n \text{ user vector}}$$

- Loss function $\mathcal{L}(\cdot)$

$$\begin{aligned} (U^t, V) &= \operatorname{argmin}_{U^t, V} \mathcal{L} \left(Z^t, V^\top U^t \right) \\ &= \left\| Z^t - V^\top U^t \right\|_2^2 \end{aligned}$$

Matrix Factorization for Modeling Users' Content Interests II

⊙ Probabilistic form

$$p(Z^t \mid U^t, V, \sigma^2) = \prod_{i=1}^n \prod_{j=1}^{p+n} \mathcal{N}(z_{ij}^t \mid v_j^\top u_i^t, \sigma^2),$$

$$p(U \mid \sigma_u^2) = \prod_{i=1}^n \mathcal{N}(u_i^t \mid 0, \sigma_u^2),$$

$$p(V \mid \sigma_v^2) = \prod_{j=1}^{p+n} \mathcal{N}(v_j \mid 0, \sigma_v^2).$$

Matrix Factorization for Modeling Users' Content Interests III

$$\begin{aligned}
 (U^t, V) = \underset{U^t, V}{\operatorname{argmin}} & \sum_{i=1}^n \sum_{j=1}^{p+n} \left\| z_{ij}^t - v_j^\top u_i^t \right\|_2^2 + \lambda_U \sum_{i=1}^n \|u_i^t\|_2^2 \\
 & + \lambda_V \sum_{j=1}^{p+n} \|v_j\|_2^2,
 \end{aligned}$$

σ^2 / σ_V^2 (under λ_V) σ^2 / σ_U^2 (under λ_U)

◎ 扩展了基础的矩阵分解模型

- 加入非线性因素;
- 考虑用户潜在因子的状态依赖关系;
- 将矩阵分解框架应用于文本建模任务中。

Matrix Factorization for Modeling Users' Content Interests IV

$g(\cdot)$ encodes the neural network parameterized by Θ

$$\begin{aligned}
 [\{U\}^{t=1:\tau}, V, \Theta] = \operatorname{argmin}_{\{U\}^{t=1:\tau}, V, \Theta} & \sum_{i=1}^n \sum_{j=1}^{p+n} \sum_{t=1}^{\tau} \left\| z_{ij}^t - g\left(v_j^\top u_i^t; \Theta\right) \right\|_2^2 \\
 & + \lambda_U \sum_{i=1}^n \sum_{t=1}^{\tau} \left\| u_i^t \right\|_2^2 + \lambda_V \sum_{j=1}^{p+n} \|v_j\|_2^2 \\
 & \underbrace{u_i^t = f(u_i^{t-1})}_{\text{blue arrow}}
 \end{aligned}$$

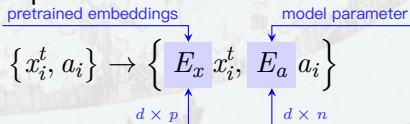
Neural Network Architecture I

$g(\cdot)$ encodes the neural network parameterized by Θ

$$\begin{aligned}
 [\{U\}^{t=1:\tau}, V, \Theta] = & \operatorname{argmin}_{\{U\}^{t=1:\tau}, V, \Theta} \sum_{i=1}^n \sum_{j=1}^{p+n} \sum_{t=1}^{\tau} \left\| z_{ij}^t - g \left(v_j^\top u_i^t; \Theta \right) \right\|_2^2 \\
 & + \lambda_U \sum_{i=1}^n \sum_{t=1}^{\tau} \left\| u_i^t \right\|_2^2 + \lambda_V \sum_{j=1}^{p+n} \|v_j\|_2^2 \\
 & \text{with } u_i^t = f(u_i^{t-1})
 \end{aligned}$$

Neural Network Architecture II

Embedding the Input Data



Estimating a Nonlinear Hidden State for Each User

$$\ell_i^t = \sigma_1 \left(W_\ell \cdot [E_x x_i^t; E_a a_i] \right)$$

Annotations:

- σ_1 is labeled as **ReLU**.
- W_ℓ is associated with dimensions $d \times 2d$.

Neural Network Architecture III

- Incorporating Dynamics by Combining a User's Current and Previous Hidden States

$$u_i^t = \underset{\text{softmax}}{\sigma_2} \left(\underset{K \times d}{W_u \ell_i^t} + \underset{K \times K}{W_r u_i^{t-1}} \right)$$

$$u_i^t = \underset{\text{smoothing hyperparameter}}{\alpha} \cdot [\sigma_2 (W_u \ell_i^t + W_r u_i^{t-1})] + (1 - \alpha) \cdot u_i^{t-1}$$

Neural Network Architecture IV

- Combining the User and Content Factors (Encoder-decoder architecture)

$$r_i^t = V^\top u_i^t$$

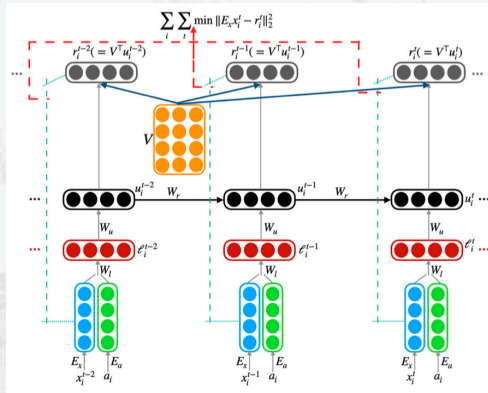
\uparrow
 $K \times d$

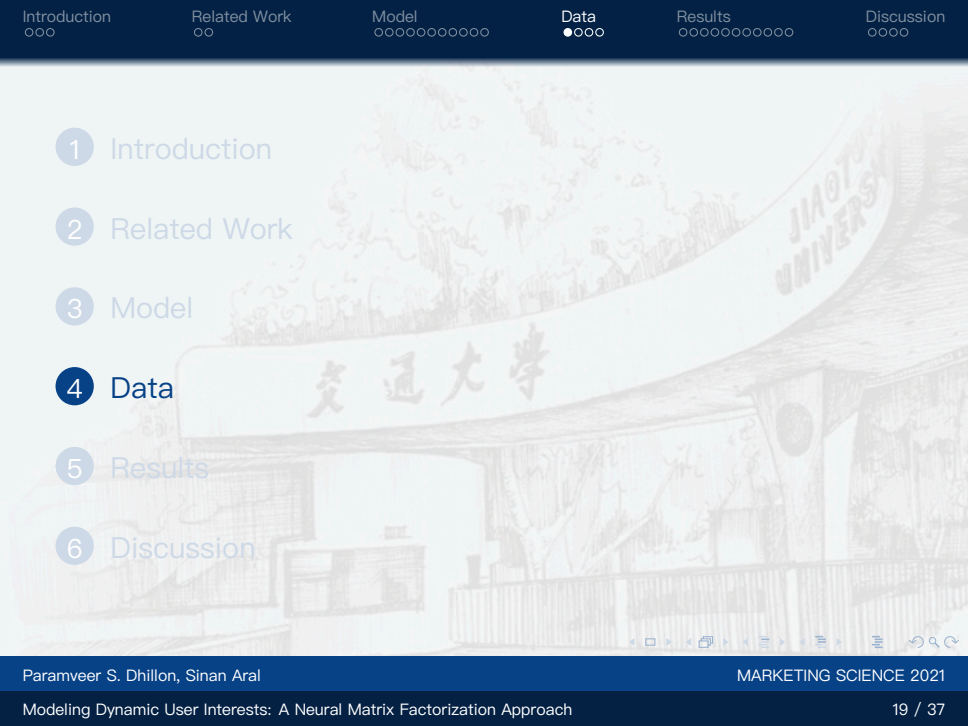
- Minimizing the Loss Function

$$[\{U\}^{t=1:\tau}, V] = \operatorname{argmin}_{U^{t=1:\tau}, V} \sum_{i=1}^n \sum_{t=1}^{\tau} \|E_x x_i^t - r_i^t\|_2^2$$

- Hyperparameters: $K = 30, \alpha = 0.5, d = 300$, random uniform initialization.

Neural Network Architecture V



- 
- A faint, artistic sketch of the Jiaotong University entrance serves as the background. It features a large, curved archway with the university's name in Chinese characters '交通大學' and English 'JIAOTONG UNIVERSITY'. There are trees and a building visible through the archway.
- 1 Introduction
 - 2 Related Work
 - 3 Model
 - 4 Data**
 - 5 Results
 - 6 Discussion

Data I

- ◎ 《波士顿环球报》2014 年 2 月 1 日-2019 年 5 月 13 日的用户点击流数据;
- ◎ 阅读的文章、阅读文章的时长、订阅状态、访问者的人口统计数据: 如地区代码、邮政编码、设备类型 (手机或台式机)、操作系统和国家;
- ◎ Week level analysis;
- ◎ Visit the website at least five different times;
- ◎ 500000 unique visitors over 276 weeks, 5610008 non-zero person-week observations;
- ◎ Use headlines, 135861569 tokens, 85228 unique tokens;

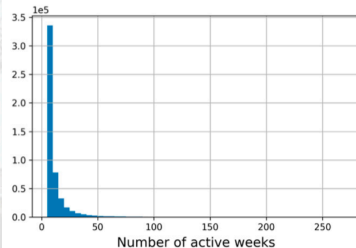
Data II

Table 1. Summary Statistics of the Visitation and Reading Behavior of the Visitors to the Globe Website

| | Min. | Median | Mean | Max. |
|------------------------|------|--------|-------|-------|
| Visits per week | 1 | 1 | 1.64 | 626 |
| News articles per week | 1 | 1 | 3.83 | 1,400 |
| Number of active weeks | 5 | 8 | 12.40 | 264 |

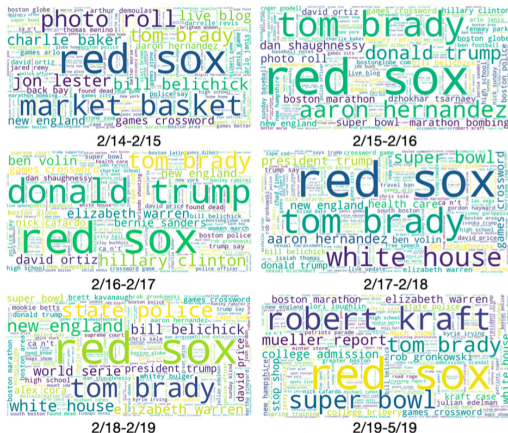
Note. Our data set consists of only those users who were active in at least five different weeks during our observation period.

Figure 1. (Color online) Frequency Distribution of User Activity



Data III

Figure 2. (Color online) Plot Showing the Prevalence of Words in the News Stories Consumed by Users in Each 52-Week (One Year) Period Starting February 2014



Note. Larger font size indicates the higher prevalence of those terms in users' consumption patterns.

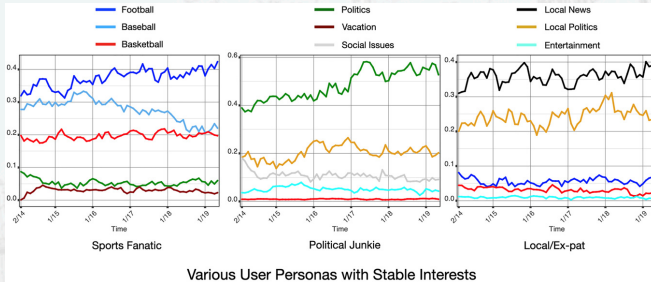
- 1 Introduction
- 2 Related Work
- 3 Model
- 4 Data
- 5 Results**
- 6 Discussion

Visualizing Trajectories of User Interests U^t

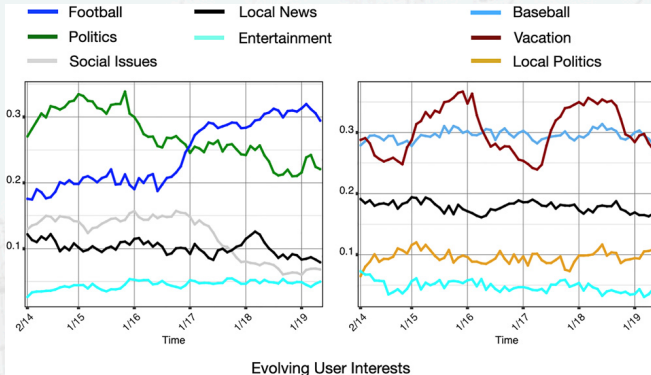
Compute the nearest neighbors of each row of V from the word embedding matrix E_x .

| Basketball | Baseball | Local News | Entertainment | Football |
|--|---|--|---|--|
| kevin nets nba ainge celtics stevens durant | sox red white yankees mariners bullpen lineup pitcher | city-hall plaza walsh mayor newton resort council | theatre art mfa ballet stage ticket museum | patriots nfl deflategate super-bowl belichik parade draft victory |
| Social Issues | Crime | Vacation | Politics | Local Politics |
| mass marijuana black race storm pike law voters | police dorchester officer charged arrested shot man killed | cape cod beaches island white mountains restaurant | clinton debate sanders poll fbi e-mails gop democratic | baker state governor charlie judge race logan policy |

Stable Interests



Evolving Interests



Evolving User Interests

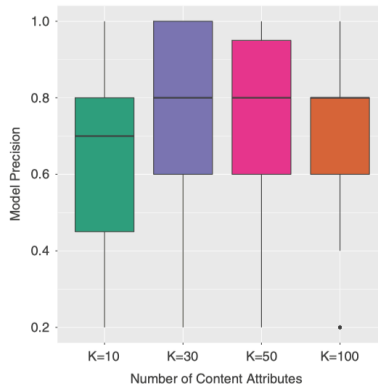
Crowdsourced Evaluation of Content Attributes I

- ◎ Word intrusion task
- ◎ $K = \{10, 30, 50, 100\}$

$$\text{Mean Precision}_k = \sum_{s=1}^S \frac{1(i_{k,s} = w_k)}{S}$$

Crowdsourced Evaluation of Content Attributes II

Figure 7. (Color online) Crowdsourced Mean Model Precision for Different Number of Attributes



Evaluating the Predictive Quality

Mean precision at K (MP@K)

$$\text{Mean Precision} = \sum_{i=1}^n \frac{1 \left(\text{NN}_1 \left(r_i^{\tau-a} \right) = c_i^{\tau} \right)}{n}$$

$t = \tau - a$ representation
 $t = \tau$ representation

nearest neighbor function

Real-valued similarity score $s(r_i^{\tau-a}, c_i^{\tau})$.

Evaluating the Predictive Quality

Table 2. Results on the Task of Retrieving the Nearest Neighbor, That Is, MP@1

| Method | $a = 1$ | $a = 2$ | $a = 3$ |
|------------------------------|----------------|----------------|----------------|
| | Mean precision | Mean precision | Mean precision |
| Weighted average of sections | 3.8 | 2.2 | 1.4 |
| LDA | 10.4 | 7.8 | 6.4 |
| LDA-GPDH | 12.2 | 10.7 | 8.7 |
| DTM | 14.9 | 12.6 | 10.9 |
| Our approach | 17.1 | 15.6 | 13.2 |

Notes. (1) Mean precision represents the fraction of users whose nearest neighbor was retrieved correctly. Please refer to Equation (11). (2) Precision numbers are multiplied by 100 to standardize them. (3) Table shows training set accuracy. (4) Model hyperparameters were tuned on the validation data set. The models are estimated on data up until a previous time periods. The prediction is always made on content consumption in the final τ^{th} period.

Table 3. Results Showing Cosine Similarity Between Embeddings of Users and the Content They Consumed

| Method | $a = 1$ | $a = 2$ | $a = 3$ |
|------------------------------|---------------------------------|---------------------------------|---------------------------------|
| | Similarity ($\mu \pm \sigma$) | Similarity ($\mu \pm \sigma$) | Similarity ($\mu \pm \sigma$) |
| Weighted average of sections | 42.9 \pm 10.6 | 40.1 \pm 9.4 | 38.7 \pm 9.2 |
| LDA | 55.4 \pm 5.1 | 52.9 \pm 4.6 | 50.1 \pm 5.4 |
| DTM | 64.6 \pm 2.1 | 61.2 \pm 2.9 | 58.6 \pm 4.3 |
| LDA-GPDH | 62.8 \pm 3.0 | 61.0 \pm 2.4 | 59.9 \pm 4.0 |
| Our approach | 71.3 \pm 3.3 | 69.4 \pm 3.9 | 67.0 \pm 3.6 |

Notes. (1) Similarity represents the cosine similarity $\mu_{u|c}^{\text{th}}$. (2) Similarity numbers are multiplied by 100 to standardize them. (3) Table shows training set accuracy. (4) Model hyperparameters were tuned on the validation data set. The models are estimated on data up until a previous time periods. The prediction is always made on content consumption in the final τ^{th} period.

Robustness Tests

Table 4. Table Showing the Impact of Hyperparameter Choice on the Validation Set Accuracy

| Hyperparameters | Mean nearest neighbor precision | | | | |
|-----------------|---------------------------------|-----------------|-----------------|-----------------|-----------------|
| | $\alpha = 0.10$ | $\alpha = 0.25$ | $\alpha = 0.50$ | $\alpha = 0.75$ | $\alpha = 0.90$ |
| K = 10 | 12.9 | 14.1 | 15.2 | 15.0 | 13.6 |
| K = 30 | 12.1 | 15.8 | 18.4 | 16.6 | 14.4 |
| K = 50 | 11.2 | 15.3 | 17.7 | 15.2 | 14.1 |
| K = 100 | 13.4 | 16.2 | 17.5 | 16.8 | 14.5 |

Notes. (1) Mean precision represents the fraction of users whose nearest neighbor was predicted correctly. Please refer to Equation (11). (2) Precision numbers are multiplied by 100 to standardize them.

Ablation Analysis

Table 5. Table Showing the Relative Contribution of Various Components of Our Model

| Ablation | Mean nearest neighbor precision | | |
|--------------------------|---------------------------------|-------|-------|
| | a = 1 | a = 2 | a = 3 |
| No nonlinearities | 11.7 | 8.6 | 6.5 |
| No time dynamics | 13.1 | 10.3 | 8.1 |
| No exponential smoothing | 15.7 | 12.9 | 10.8 |
| Full model (Table 2) | 17.1 | 15.6 | 13.2 |

Notes. (1) Mean precision represents the fraction of users whose nearest neighbor was predicted correctly. Please refer to Equation (11). (2) Precision numbers are multiplied by 100 to standardize them. Training set accuracy is reported. The models are estimated on data up until a previous time periods.

Scalability, Transferability, and Cold-Starting New Users.

- ⊙ Use transfer learning to learn representations for users with consumption traces x_i^t .
- ⊙ Freezing all the estimated parameters: $\{W_\ell, W_u, W_r, V, E_x\}$.
- ⊙ Use the observable user demographics.

- 1 Introduction
- 2 Related Work
- 3 Model
- 4 Data
- 5 Results
- 6 Discussion**

Managerial Implications

- ⊙ Generating User Profiles
- ⊙ Content Categorization and Recommendation

Conclusion and Limitations

- ◎ Proposed a neural matrix factorization modeling approach to extract nonlinear patterns from text data to infer customers' evolving interests.
- ◎ Model other types of high-dimensional consumption data(e.g., images, videos, and audio).
- ◎ Limitations
 - Model only one kind of customer digital footprints—news consumption(?);
 - Deep learning and neural net models in marketing research is still an under explored area of study;
 - Only model the demand side and assume that consumers' consumption patterns are driven only by their consumption in previous periods.

Thanks!