

Transport optimal appliqué à la recommandation d'œuvres

RYAN LAHFA

Table des matières

1	Position du problème	1
1.1	Recommandation d'entrées	1
1.2	Jeu de données : Mangaki	1
2	Modèle de comparaison : 20-KNN	2
2.1	Introduction	2
2.2	Choix du paramètre k , de la métrique et visualisation des voisins	2
2.3	Défauts et limites du modèle	2
2.4	Objectifs du TIPE	2
2.5	État actuel de la recherche	3
3	Raffinement par le transport optimal : impact de la distance de Wasserstein	3
3.1	Propriétés de \mathcal{W}	3
3.2	Intérêt : calcul efficace et rapide \mathcal{W} , propagation de l'information visuelle dans le modèle	3
3.3	Calcul des représentations visuelles par le réseau de neurones convolutifs Illustration2Vec	3
4	Résultats	4
4.1	AUC	4
4.2	Temps de calcul	4
4.3	Analyse qualitative	4
5	Prolongements envisagables	4
	Références bibliographiques	4

Table des figures

1 Position du problème

1.1 Recommandation d'entrées

À partir d'un ensemble de préférences exprimés par des utilisateurs, l'on veut une méthode de prédire les futures préférences des utilisateurs, c'est le problème qu'on essayera de résoudre.

Précisément, on se donne une base de données représentée par une matrice $M \in \mathcal{M}_{n,m}(\{0,1\})$ dont le terme général $(m_{i,j})$ indique si l'utilisateur i a aimé l'entrée j .

À partir de M , l'on veut apprendre un classificateur capable de prédire $\widehat{m}_{i,j}$ si celui-ci n'est pas connu, avec, ou non, $p_{i,j} = \mathbb{P}(m_{i,j} = \widehat{m}_{i,j}) \in [0,1]$.

On notera aussi $r_u \in \mathbb{R}^m$ la distribution de l'utilisateur $u \in [[1,n]]$ sur les entrées.

1.2 Jeu de données : Mangaki

Le modèle sera testé sur le jeu de données fournis par le site [5] qui comporte :

- 2289 utilisateurs ;
- 12479 œuvres issus de l'animation japonaise (animes, mangas) ;
- plus de 350000 notes

2 Modèle de comparaison : 20-KNN

2.1 Introduction

Pour $k \geq 1$ (ici $k = 20$ d'où 20-KNN), le modèle des k -plus proches voisins consiste en :

- Pour chaque utilisateur i , calculer des voisins, qu'on notera $\mathcal{N}(i) \subset [[1, n]]^k$ au sens d'une métrique $d : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}_+$ opérant sur les distributions d'utilisateurs, prendre les k plus proches
- Pour classer une nouvelle entrée pour un utilisateur i , on fait « voter » les voisins de i et on prend la majorité comme prédiction, i.e.

En notant $\mathcal{N}'(i)_j = \{j \in \mathcal{N}(i) \mid m_{u,j} \text{ est connu}\}$ et m son cardinal.

$$\widehat{m}_{i,j} = \begin{cases} 1 & \text{si } 2 \sum_{k \in \mathcal{N}'(i)_j} m_{k,j} \geq m \\ 0 & \text{sinon.} \end{cases}$$

On abrégera KNN pour le modèle des k -plus proches voisins dans le reste du document et 20-KNN pour $k = 20$.

2.2 Choix du paramètre k , de la métrique et visualisation des voisins

Le choix du paramètre k peut s'effectuer par validation croisée sur le jeu de données, cette validation croisée a été effectuée préalablement et fournit que $k = 20$ donne une bonne performance relativement à la racine carrée de l'erreur moyenne au carré (RMSE).

Ensuite, pour la métrique, on utilise la similarité cosinus qui possède de bonnes performances empiriquement sur les tâches de recommandation d'après [4], dont on rappelle la définition :

$$\text{sim}(u, u') = \frac{r_u^\top r_{u'}}{\|r_u\|_2 \|r_{u'}\|_2}$$

où $\|\cdot\|_2$ est la norme ℓ_2 .

On peut aussi procéder à une visualisation des graphes de voisins :

2.3 Défauts et limites du modèle

En introduction, l'entraînement de 20-KNN dépend de la métrique d employée, si on note $\text{Supp}(u)$ pour $u \in [[1, n]]$ le support des utilisateurs, défini par :

$$\text{Supp}(u) = \{j \in [[1, m]] \mid m_{u,j} \text{ est connu}\}$$

Alors, pour u, v deux utilisateurs tels que $\text{Supp}(u) \cap \text{Supp}(v) = \emptyset$, alors : $\text{sim}(u, v) = 0$

Or, la situation dans laquelle l'utilisateur u a lu les versions mangas d'une œuvre et v a vu les versions animées de celle-ci peut se présenter, cependant la métrique n'en tient pas compte et ne peut le calculer puisqu'il s'agit d'une information propre à l'œuvre.

2.4 Objectifs du TIPE

Nous répondrons aux questions suivantes :

- Sachant qu'on dispose de l'ensemble des couvertures des œuvres, peut-t-on calculer une métrique qui tient compte de l'information visuelle de ces couvertures et des similarités entre les distributions d'utilisateurs ?
- En la remplaçant par la similarité cosinus, obtient-t-on une meilleure performance au sens d'une métrique d'erreur ?
- Est-ce qu'on constate des transferts d'information pertinents et intéressants tels que : la saison i d'une œuvre vers la saison $i + j$ de la même œuvre, du format manga vers le format anime ou vice versa ?

2.5 État actuel de la recherche

À notre connaissance, la littérature ne mentionne pas beaucoup de travaux qui cherchent à intégrer des métadonnées visuelles dans un système de recommandation afin d'en améliorer sa qualité et son interprétabilité, on notera [2].

3 Raffinement par le transport optimal : impact de la distance de Wasserstein

Le transport optimal est un domaine qui est de plus en plus appliqué notamment grâce à [1] qui a permis le calcul effectif et approximatif des objets de façon tractable.

Si l'on dispose de $r \in \mathbb{R}^d, c \in \mathbb{R}^d$ deux distributions de probabilités discrètes, en posant $U(r, c) = \{M \in \mathcal{M}_{d,d}(\mathbb{R}_+) \mid M \mathbb{1}_d = r \text{ et } M^\top \mathbb{1}_d = c\}$, l'ensemble des probabilités jointes sur r et c , on définit la distance de Wasserstein comme étant :

$$\mathcal{W}(r, c) = \min_{\gamma \in U(r, c)} (\gamma \mid C)_F$$

où C est une matrice exprimant le coût de transporter de la masse de r_i vers c_j et $(\cdot \mid \cdot)_F$ est le produit scalaire de Frobenius.

3.1 Propriétés de \mathcal{W}

La distance de Wasserstein est une métrique.

3.2 Intérêt : calcul efficace et rapide \mathcal{W} , propagation de l'information visuelle dans le modèle

Par l'algorithme de Sinkhorn, présenté en détails dans [1], il est possible de calculer une approximation de \mathcal{W} , pour $\varepsilon > 0$, un paramètre de régularisation entropique :

$$\mathcal{W}_\varepsilon(r, c) = \min_{\gamma \in U(r, c)} (\gamma \mid C)_F + \varepsilon \Omega(\gamma)$$

On prouve aussi dans [1] que $(r, c) \mapsto \mathbb{1}_{r \neq c} \mathcal{W}_\varepsilon(r, c)$ est une distance.

3.3 Calcul des représentations visuelles par le réseau de neurones convolutifs Illustration2Vec

En utilisant les travaux de [3] et [6], on peut calculer des représentations parcimonieuses de couvertures $(p_i)_i \in (\mathbb{R}^{512})^m$ qui permettent de poser la matrice de coût comme étant :

$$C = (\|p_i - p_j\|_2^2)_{(i,j) \in [[1,m]]^2}$$

Donc, la matrice de coût représente la similarité visuelle entre deux couvertures.

4 Résultats

4.1 AUC

4.2 Temps de calcul

4.3 Analyse qualitative

5 Prolongements envisagables

Références bibliographiques

- [1] Cuturi, M. 2013. Sinkhorn distances : Lightspeed computation of optimal transport. *Advances in neural information processing systems* (2013), 2292-2300.
- [2] Messina, P., Dominguez, V., Parra, D., Trattner, C. et Soto, A. 2019. Content-based artwork recommendation : integrating painting metadata with neural and manually-engineered visual features. *User Modeling and User-Adapted Interaction*. 29, 2 (2019), 251-290.
- [3] Saito, M. et Matsui, Y. 2015. Illustration2vec : a semantic vector representation of illustrations. *SIGGRAPH Asia 2015 Technical Briefs* (2015), 5.
- [4] Sarwar, B.M., Karypis, G., Konstan, J.A., Riedl, J. et others 2001. Item-based collaborative filtering recommendation algorithms. *Www*. 1, (2001), 285-295.
- [5] Vie, J.-J., Lailly, C. et Pichereau, S. 2015. Mangaki : an anime/manga recommender system with fast preference elicitation. *Tech. Rep.* (2015).
- [6] Vie, J.-J., Yger, F., Lahfa, R., Clement, B., Cocchi, K., Chalumeau, T. et Kashima, H. 2017. Using posters to recommend anime and mangas in a cold-start scenario. *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)* (2017), 21-26.