

THEMATIC RECOMMENDATIONS ON KNOWLEDGE GRAPHS USING MULTILAYER NETWORKS*

MARIANO BEGUERISSE-DÍAZ^{†‡}, DIMITRIOS KORKINOF^{†‡}, AND TILL HOFFMANN

Abstract.

We present a framework to generate and evaluate thematic recommendations based on multilayer network representations of knowledge graphs (KGs). In this representation, each layer encodes a different type of relationship in the KG, and directed interlayer couplings connect the same entity in different roles. The relative importance of different types of connections is captured by an intuitive salience matrix that can be estimated from data, tuned to incorporate domain knowledge, address different use cases, or respect business logic.

We apply an adaptation of the personalised PageRank algorithm to multilayer models of KGs to generate item-item recommendations. These recommendations reflect the knowledge we hold about the content and are suitable for thematic and/or cold-start recommendation settings. Evaluating thematic recommendations from user data presents unique challenges that we address by developing a method to evaluate recommendations relying on user-item ratings, yet respecting their thematic nature. We also show that the salience matrix can be estimated from user data. We demonstrate the utility of our methods by significantly improving consumption metrics in an AB test where collaborative filtering delivered subpar performance. We also apply our approach to movie recommendation using publicly-available data to ensure the reproducibility of our results. We demonstrate that our approach outperforms existing thematic recommendation methods and is even competitive with collaborative filtering approaches.

Key words. Knowledge Graphs, Multilayer Networks, Thematic Evaluation, Recommendation Systems

AMS subject classifications. 05C21, 05C82, 60J20, 91D30

1. Introduction. Knowledge graphs (KGs) encode semantic relationships between large collections of items [19, 27]. In recent years, there has been a surge of interest in KGs in both academic and industrial settings. Much of this interest is due to a shift from information to knowledge retrieval, and initiatives like the Google knowledge vault [15] or Amazon’s product graph [16]. Typical applications of KGs include reasoning [4, 9], search [52], and recommendations [54].

Recommendation systems fall into three broad categories: (i) collaborative filtering (CF) approaches [44], formulated as a missing data problem to predict the likely interaction between a user and an item, (ii) reinforcement learning (RL) approaches [55, 35] to continuously refine recommendations in an online setting, and (iii) knowledge-based approaches, which exploit the information we hold about items to generate recommendations [44, 10]. Graph-based approaches often fall into the latter category, and knowledge is captured by relationships among items and other entities (e.g. tags or users) [27]. Some recommender systems fuse multiple approaches to further improve recommendations: for example, CF and KGs [36] or KGs with RL [51].

Collaborative filtering systems predict which items users are most likely to interact with based on their past interactions with other items [31, 40]. The success of CF systems has made them a mainstay in research and practice of recommender systems. However, there are situations where user-item interactions alone are not enough to achieve the desired performance [10], or sufficient interaction data may not be available

*Submitted to the editors May 13, 2021.

[†]Spotify Ltd, UK. (marianob@spotify.com, dkorkinof@spotify.com)

[‡]These authors contributed equally to this work

(e.g. a cold-start recommendation setting) [23]. Several recommendation systems thus combine collaborative filtering approaches with annotations and KGs to enrich their recommendations [47, 50].

A common approach to recommendation systems based on KGs is to generate node embeddings using graph networks or matrix factorisations [2, 17, 41, 43, 46, 53, 54]. Although embedding methods can achieve impressive performance, their recommendations can be difficult to explain. This has motivated researchers to work on methods to understand and explain recommendations [2]. Furthermore, the geometry of embedding spaces can be complex [7], which presents difficulties in tasks such as proximity search and multi-seeding, i.e. using multiple items to provide a recommendation context. The lack of interpretability furthermore poses challenges for assessing fairness and bias in complex machine learning systems [11, 39]. From a modelling perspective, one concern that arises in some embedding methods is that they may not fully exploit the richness of the structure of the KG [41], or that they do not adequately distinguish between connections of different type in the KG [54]. Some methods do use the graph structure to generate recommendations, e.g. using random walks and diffusion processes [50, 8, 37, 20]; however, these works do not often make a distinction between different types of connections, which can have an important effect on the quality of the recommendations.

In this work we propose a principled method to generate *thematic* item-item recommendations using random walks on multilayer network representations of KGs. Broadly, we say recommendations are thematic if they rely on the intrinsic properties of the items, which can be encoded as connections and nodes in a KG. Our method enables thematic exploration of items in the KG with or without user data, making it applicable to cold-start recommendation problems. In addition, our method allows us to exploit the full structure of the KG, accounting for the distinction between different types of connections.

This paper is structured as follows: in Sec. 2, we show how to represent a KG as a multilayer network where each layer represents one connection type; nodes that represent the same entity in different layers are connected to each other by directed, weighted interlayer couplings. Then, in Sec. 3 we construct the rate matrix of a random walk on the multilayer network, incorporating a set of parameters that encode the salience of each connection type. These saliences put different connection types on an equal footing, facilitating direct comparison between connections. In Sec. 4, we use an adaptation of personalised PageRank to produce recommendations using arbitrary sets of nodes as seeds. We introduce a method to evaluate thematic recommendations based on user data in Sec. 5; this evaluation differs from traditional CF evaluation approaches because it assesses recommendations generated using the connections in the KG only. In Sec. 6, we showcase two applications of our method: in Sec. 6.1, we show the results of an AB test for thematic music recommendations on a cohort of 100k users, and in Sec. 6.2, we apply our method to a KG of the film industry, which we evaluate offline using the method from Sec. 5. We show that we can learn the value of the salience parameters by optimising this evaluation score. These saliences can be interpreted as the relative importance of the different types of connections in the KG, and can be used to explain the output of the system. We compare our approach with alternative methods and show that our recommendations can outperform other thematic recommendation systems and be competitive with CF methods. Finally, in Sec. 7, we summarise our methods and results, and indicate potentially interesting avenues for future research and applications.

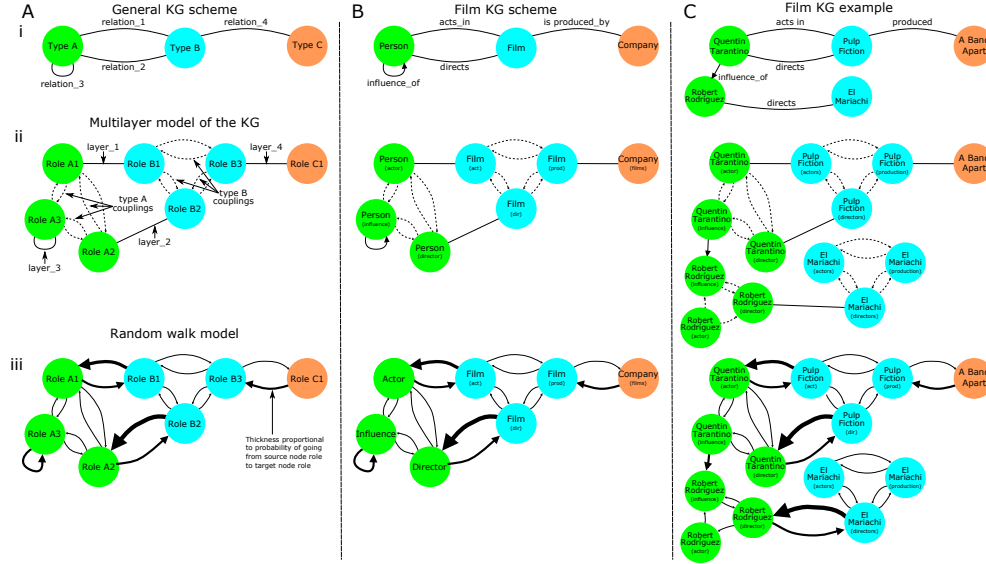


Fig. 1: **A:** (i) A general schematic of a KG with three types of entities and four types of relations. (ii) Multilayer network representation of the KG. Each layer corresponds to a type of relation in the KG (solid lines). Unipartite layers model interactions among nodes representing the same role, and bipartite layers model interactions between two different roles. Nodes that represent the same entity in different layers are coupled to each other (dashed lines). (iii) Illustration of a random walk model on the multilayer representation of the KG. In this diagram, the width of each directed connection indicates the proclivity of a random walker to go from a node in the source role to another in the target role. **B:** a schematic of a film KG with three types of entities (person, film, company), and four types of relations (acting, directing, influencing and producing). In the multilayer representation of the KG, there are three roles for each person entity (director, actor, influence), three roles for films (one for each of its interactions with actors, directors and companies), and one role for companies. **C:** a specific example of a film KG, its multilayer representation, and random walk model.

2. Network models of knowledge graphs. Knowledge graphs (KGs) are typically *multigraphs*: entities, represented by nodes, are connected by different types of relationships [14]. Figure 1 shows an example of a schematic of a KG about the film industry. This KG comprises three types of entities: person, film, and company. Companies produce films; a person can contribute to a movie by either directing it or acting in it, and influence other people’s work. These types of relations encode distinct information and are thus neither interchangeable nor directly comparable. For example, naively evaluating the centrality of nodes in a KG by treating all relationships in the same manner would discard important information, and may lead to suboptimal results. A recommendation algorithm using KGs should account for the different types of relationships.

To make the best use of all the information encoded by the relationships in the KG, yet develop a tractable thematic recommendation system, we adopt a multilayer network approach [29]: each relationship type is represented by one layer, and the same entity in the KG can appear in different layers, fulfilling different roles. Hence-

forth, we will use the term *role* to refer to how an entity in the multilayer network relates to others (e.g. a person acting in or directing a movie) and reserve the term entity *type* to refer to what a node is (e.g. a person, company or movie).

The knowledge graph in Fig. 1 encodes information about a set \mathcal{P} of person entities, \mathcal{C} of company entities, and \mathcal{F} of film entities with cardinality P , C , and F , respectively. Person and company entities play one of the following roles:

- people acting in movies are represented by **actor** roles (**p-act**),
- people directing movies are represented by **director** roles (**p-dir**),
- people influencing others are represented by **influencer** roles (**p-inf**),
- and companies producing movies are represented by **production company** roles (**c-prod**).

The passive involvement of movies is represented by the following complementary roles:

- **film-acting** for films connected to actors (**f-act**),
- **film-directing** for films connected to directors (**f-dir**),
- and **film-producing** for films connected to production companies (**f-prod**).

Nodes in the multilayer network that represent the same entity, such as Quentin Tarantino being both the director of and an actor in the 1994 film Pulp Fiction, are connected to one another by interlayer couplings. Figure 1B provides an illustration of the entities and relationships outlined above.

More formally, we consider a multilayer network of the film KG in Fig. 1 with four layers, one for each relationship type, and $N = 3F + 3P + C$ nodes across all layers: one node for each entity in all roles in which it *could* be active (e.g. an actor who does not direct still has a director node). The supra-adjacency matrix $A \in \mathbb{R}^{N \times N}$, i.e. a flattened representation of the multilayer network [29], is defined by

$$(2.1) \quad A = B + \Lambda,$$

where B is a block matrix representing the relationships of the knowledge graph. Each block $B^{(l)}$ in B is the adjacency matrix of the layer encoding relationship type l (i.e. the *intralayer* connections in a multilayer network). The size of $B^{(l)}$ depends on the entity types that participate in l . For example, the **directs** block represents a bipartite graph and has P rows and F columns. The weight B_{ij} of the connection between two nodes $i, j \in \{1, \dots, N\}$ encodes the a priori importance of the relationship, such as the fraction of screen time an actor has in a movie. If the nodes are not connected, then $B_{ij} = 0$.

The block matrix Λ captures the directed, weighted *interlayer* couplings between all nodes that represent the same entity (illustrated by dashed lines in Fig. 1). The directed interlayer coupling from j to i exists if the two nodes represent the same entity in different roles, and the corresponding weight is equal to the weighted degree of node i within its layer. Consequently, nodes that are not active in a given role have no incoming interlayer connections, but do have outgoing ones. For example, a connection from Tarantino’s actor node to Tarantino’s director node exists because they represent the same person *and* he has directed movies. In contrast, Robert Rodríguez has no acting credits (in this example), thus his actor node has no incoming interlayer connections, as shown in Fig. 1C. Each block of Λ is a square matrix $\Lambda^{(r)}$ whose main diagonal contains the weighted degree of all nodes with role r , and has zeroes everywhere else. The size of each block depends on the entity type that assumes the role. For example, $\Lambda^{(\text{p-dir})}$ has size $P \times P$, and $\Lambda_{ij}^{(\text{p-dir})} = k_i^{(\text{p-dir})} \delta_{ij}$ for $i, j \in$

$\{1, \dots, P\}$ where $k_i^{(p\text{-dir})}$ is the weighted degree of person i as director and δ_{ij} is the Kronecker delta.

Finally, the block structure of the supra-adjacency matrix of the multilayer network is

$$(2.2) \quad A = \begin{matrix} & \text{f-act} & \text{f-dir} & \text{f-prod} & \text{p-act} & \text{p-dir} & \text{p-infl} & \text{c-prod} \\ \text{f-act} & & & & & & & \\ \text{f-dir} & \Lambda^{(\text{f-dir})} & & & & & & \\ \text{f-prod} & \Lambda^{(\text{f-prod})} & \Lambda^{(\text{f-prod})} & & & & & \\ \text{p-act} & B^{(\text{act})} & & & & & & B^{(\text{prod})\top} \\ \text{p-dir} & & B^{(\text{dir})} & & & & & \\ \text{p-infl} & & & & \Lambda^{(\text{p-dir})} & \Lambda^{(\text{p-act})} & \Lambda^{(\text{p-act})} & \\ \text{c-prod} & & & B^{(\text{prod})} & \Lambda^{(\text{p-infl})} & \Lambda^{(\text{p-infl})} & B^{(\text{infl})} & \end{matrix}.$$

This matrix encodes the KG in Fig. 1 *without loss of any information*. For example, if Quentin Tarantino and Robert Rodríguez are $i, k \in \mathcal{P}$ respectively, and Pulp Fiction is $j \in \mathcal{F}$, the acting connection between Tarantino and Pulp Fiction in Eq. (2.2) is found in $A_{j, 3F+i}$ and $A_{3F+i, j}$ (via $B_{ij}^{(\text{act})}$); the directing connection is in $A_{F+j, 3F+P+i}$ and $A_{3F+P+i, F+j}$ (via $B^{(\text{dir})}$); Tarantino's influence on Rodríguez is only in entry $A_{3F+2P+k, 3F+2P+i}$ (via $B^{(\text{infl})_{ik}}$) because influence is a directed layer. Finally, the fact that Tarantino the actor, the director, and the influence are the same person is encoded in six elements of A (via $\Lambda^{(\text{p-dir})_{ii}}$, $\Lambda^{(\text{p-act})_{ii}}$ and $\Lambda^{(\text{p-infl})_{ii}}$). Note that an entity that is not active in a specific role still has a node for it; this node has no incoming connections, and the only outgoing connections are interlayer couplings to itself in roles in which the entity is active. An assumption we make throughout this work is that all entities are active in at least one layer (i.e. there are no isolated entities in the KG).

We can reduce the size of the multilayer model by merging roles. For example, if we merge all movie roles into a single one, the supra adjacency of the multilayer network would be

$$(2.3) \quad A' = \begin{matrix} & \text{film} & \text{p-act} & \text{p-dir} & \text{p-infl} & \text{c-prod} \\ \text{film} & & & & & \\ \text{p-act} & B^{(\text{act})} & & & & \\ \text{p-dir} & B^{(\text{dir})} & \Lambda^{(\text{p-dir})} & & \Lambda^{(\text{p-act})} & \\ \text{p-infl} & & \Lambda^{(\text{p-infl})} & \Lambda^{(\text{p-infl})} & B^{(\text{infl})} & \\ \text{c-prod} & B^{(\text{prod})} & & & & \end{matrix},$$

with $F + 3P + C$ nodes. In A' , all the interactions with movies are in the first block row/column, and that there are no interlayer couplings Λ between movie nodes. The most aggressive merging strategy would create a single role for each entity (i.e. one node per person, one per film, and so on), and would have $F + P + C$ nodes (compared with $3F + 3P + C$ nodes in the full model):

$$(2.4) \quad A'' = \begin{matrix} & \text{film} & \text{person} & \text{company} \\ \text{film} & & & \\ \text{person} & B^{(\text{act})} + B^{(\text{dir})} & B^{(\text{act})\top} + B^{(\text{dir})\top} & B^{(\text{prod})\top} \\ \text{company} & B^{(\text{prod})} & B^{(\text{infl})} & \end{matrix}.$$

Merging roles can have computational and modelling benefits in certain applications. However, because connections are not generally comparable (e.g. acting is not the same as directing), roles need to be merged with care. Furthermore, it may not be

each block, i.e. $\alpha^{(r_1, r_2)}_{ij} = c^{(r_1, r_2)} \in \mathbb{R}^+$ for all i and j , and roles r_1 and r_2 , although it is possible to encode saliences for each node, or even individual connections in the most general case. Note that the connections in the KG may be dimensional, so the units of the corresponding $\alpha^{(r_1, r_2)}$ matrix should be the inverse of the connections'. For example, if the weight of connections from **actor** to **film-acting** nodes encodes minutes on screen, the units of the corresponding entries of $\alpha^{(\text{f-acting}, \text{act})}$ are min^{-1} .

The first step to construct the transition matrix of a random walk is to combine both the salience of different relationship types and the knowledge encoded in the graph. We define the matrix

$$(3.2) \quad \tilde{T} = A \circ \alpha,$$

where \circ denotes the element-wise product. Because we combine the a priori connection strength and salience, the entries of \tilde{T} are non-dimensional and directly comparable within each column, which means that we can now normalise its columns to construct the rate matrix of a random walk. The probability that a walker transitions from node j to node i is

$$(3.3) \quad T_{ij} = \frac{\tilde{T}_{ij}}{\sum_{i=1}^N \tilde{T}_{ij}},$$

which defines a matrix T equal to \tilde{T} after column-normalisation to guarantee that the density of walkers is conserved; i.e. T is column-stochastic. These probabilities depend only on the KG data and the salience matrix α in Eq. 3.1. Therefore, we can exercise some control on the dynamics of a random walk by tuning α ; this can be useful in contexts where different aspects of the data become more or less important (e.g. emphasising directing connections over acting in one context, and vice versa in another). See Refs. [12, 48] for similar ways to parametrise random walks on multilayer networks.

To rank nodes in the multilayer model of the KG, we adopt the widely-used personalised PageRank model of Ref. [38]: a walker follows a link with probability $1 - \rho \in [0, 1]$ and ‘‘teleports’’ to another node with probability ρ . The probability to transition to a particular node via teleportation is encoded in the $N \times 1$ vector v , where $v_i \geq 0$ is the probability that node i receives a teleported walker, and $\sum_i v_i = 1$. Let $x_i(t)$ denote the probability that a walker is present on node i at time t ; this vector evolves according to the transition rule

$$(3.4) \quad x(t+1) = (1 - \rho)Tx(t) + \rho v.$$

The PageRank vector is the steady state of the distribution of walkers π , which occurs when $x(t+1) = x(t) = \pi(v)$. The value of each element $\pi_i(v)$ is the fraction of time that a random walker spends on node i in an infinitely long random walk, given a teleportation vector v . Substituting $\pi(v)$ into Eq. 3.4, we obtain the linear system [13]

$$(3.5) \quad (I - (1 - \rho)T)\pi(v) = \rho v,$$

where I is the identity matrix.

Intuitively, the teleporation vector v encodes the seeds of a walker on the multilayer network, and ρ captures the curiosity of the walker, i.e., the likelihood it will explore the network. In the limit $\rho \rightarrow 0$, we recover eigenvector centrality, and the steady state state is trivially equal to the teleportation vector when $\rho = 1$. While

the steady state solution can in principle be obtained by iterating Eq. 3.4 until convergence, the formulation in Eq. 3.5 is preferable because we can use state of the art numerical techniques to solve the linear system [13]. See Ref. [22] for a review on PageRank.

4. Thematic recommendations on multilayer networks. The fundamental assumption underpinning this work is that the connections in the KG are relevant for recommendations. Given an item or set of items of interest (i.e. the seeds), we can retrieve other items in the KG that are related to them via the graph’s connections. We can encode the recommendation context in the teleportation vector v by giving seed nodes a higher probability of receiving teleported walkers. For example, the teleportation vector appropriate for a user who has expressed interest in both Quentin Tarantino and Samuel L. Jackson will assign substantial weight to their corresponding nodes. Consequently, the steady state distribution of the random walk $\pi(v)$ in Eq. (3.4) will be biased in favour of nodes that are easily reachable from those seeds. The value of $\pi_i(v)$ can be interpreted as a contextual score of importance for node i given the seed vector v . Thus, we can obtain a ranked list of recommendations as

$$z'(v) = \text{argsort}(-\pi(v)).$$

For our hypothetical Tarantino-Jackson aficionado, Pulp Fiction is likely to have a high probability of being visited given the seed vector, and it would appear near the top the ranked list of nodes.

More generally, the teleportation vector v depends on the recommendation context, and it may be determined by implicit feedback (such as viewing behaviour), explicit feedback (such as ratings), or direct user input (such as a semantic search query). This method is thus versatile and applicable in a range of information retrieval settings. Representing the different roles of a given entity as separate nodes in a multilayer network provides unique opportunities for fine-grained recommendations. For example, we can seed recommendations with Quentin Tarantino as an actor or director depending on the preferences of the user.

While intuitive, centrality measures based on random walks, such as PageRank, suffer from localisation [33], i.e. often most of the mass of π is associated with a small number of highly-connected hubs. This behaviour is undesirable for thematic recommendations because a large number of KG connections is a result of rich knowledge about the node rather than an intrinsic measure of quality (unlike in the original formulation of PageRank, in which directed connections between websites can be seen as declarations of interest or quality). To attenuate the effect of hubs on recommendations, we can filter the ranked list of recommendations such that

$$(4.1) \quad z(v) = \left\{ i \in z'(v) : \log_{10} \frac{\pi_i(v)}{\pi_i^*} \geq \theta \right\},$$

where $\pi_i(v)$ denotes the seeded PageRank score of node i , the vector π^* denotes the unseeded PageRank score (i.e. when $v_i = \frac{1}{N}, \forall i$), and $\theta \in \mathbb{R}$ is a threshold. The filtering in Eq. 4.1 serves to eliminate candidates that are not sufficiently “thematically” related to the seeds; that is, their PageRank score did not increase enough as a result of the seeding (compared to the unseeded PageRank scores). For example, a threshold $\theta = 0$ ensures that the PageRank score of a candidate does not decrease after seeding; $\theta = 1$ keeps only nodes whose score increased by at least an order of magnitude. A negative θ tolerates some decrease in scores; this can be useful in contexts where the presence of hubs or popular nodes in the recommendations is desirable.

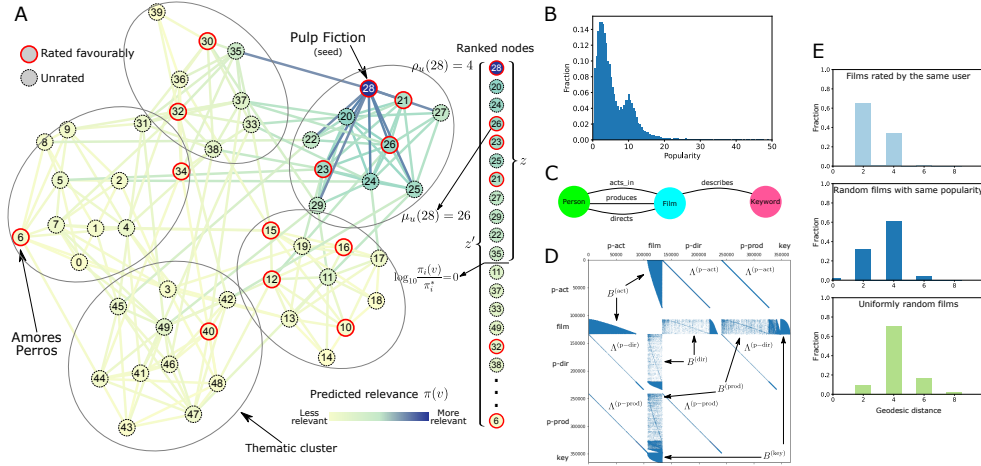


Fig. 2: **A.** An example network where connections between nodes denote a thematic relationship. A user’s favourably-rated items (\mathcal{S}_u) appear circled in red. The colour of each node represents its entry in $\pi(v)$. Seeding on Pulp Fiction (node 28) gives the ranking of the nodes z' to the right, we also indicate the set z of nodes whose PageRank increases as a result of the seeding (i.e. $\theta = 0$ in Eq. 4.1). Given the seeding with Pulp Fiction, the highest-ranked item that was also rated positively by the user is node 26 at rank 4, i.e. $\mu_u(28) = 26$ and $\rho_u(28) = 4$. **B.** Distribution of popularity of films in TMDb. **C.** Schematic of the KG we construct from TMDb with three entity types and four connection types. **D.** Sparsity pattern of the supra adjacency matrix of multilayer network from the TMDb KG. We indicate which blocks correspond to which layers and interlayer couplings. **E.** Top: distribution of the shortest geodesic distance between films rated by 1,000 Movielens users (i.e. $\text{dist}(s, s')$ where $s, s' \in \mathcal{S}_u$). Middle: shortest geodesic distance between films in 1,000 random sets of films with the same size and popularity distribution as the sets rated by the users in the top panel. Bottom: shortest geodesic distance between films in 1,000 random sets of films with the same size as the user sets, where the films were chosen uniformly at random.

5. Offline evaluation of thematic recommendations from user-item interactions. One of the best ways to evaluate the quality of a recommendation system is with an AB test, in which one cohort of users is exposed to recommendations from the system we wish to evaluate, and users from the control group receive recommendations from a system in place (or an alternative system). Although large-scale AB tests are an excellent way of evaluating the quality of recommendations, they are costly, time consuming, and risk sub-par user experiences for a proportion of users [45]. As a result AB, tests are not practical for estimating good values of the unknown parameters of the network, such as the salience matrix α or teleportation probability ρ . In addition, AB testing may not be possible without access to a platform with enough users, which is also a barrier to researchers that wish to reproduce these results.

An alternative to AB tests lies in the judicious analysis of user-item interaction data. However, using such data is challenging because consumption can be biased in favour of the recommendation system in place [21, 24], governed by personal user preference (as opposed to *thematic* similarity), and driven by popularity [3]. Never-

theless, thematic signals often do exist in user data (as we show in Sec. 6.2.1 for movie recommendation), although separating them from user preferences is non-trivial.

5.1. Extracting thematic relevance from user preference data. We extract thematic information from user-item interaction data by relying on two main assumptions: 1) the KG is by construction “thematically clustered” (i.e. the KG has communities), and 2) users show affinity towards a number of thematic clusters (i.e., users have thematic preferences such as favourite genres, actors, directors and so on), which manifests itself as favourable ratings or increased consumption of items concentrated in these clusters.

Figure 2A illustrates these assumptions in a KG, where nodes are connected if they are thematically related (e.g. share an actor). The items rated favourably by the user (in red) are distributed across the KG, but are denser in the clusters that the user has most affinity to. In a traditional CF recommendation system, we would expect all the nodes in red to be near each other in a ranking. For example the 2000 film *Amores Perros* is often compared to *Pulp Fiction* (the seed) [18]; it is likely that both films are rated similarly by several users, which a CF method would use as a signal to recommend them together. However, these two films are thematically distant and would be too far apart in the KG to be recommended together by our method.

Given a set \mathcal{M}_u of items that user u finds relevant (i.e. rated favourably or frequently consumed), we can use one item $m \in \mathcal{M}_u$ as a seed and record the rank of every other item $m' \neq m$, in the list of recommendations $z(v_m)$, where v_m is the teleportation vector with most of its mass concentrated at item m , obtained by using Eq. (4.1). We denote these ranks as $r_{mm'} = \text{rank}(m'|z(v_m))$. The higher the ranks $r_{mm'}$ are, the better recommendation list $z(v_m)$ is for this particular user. It would be tempting to use well-known scores such as the Mean Average Precision, or the Normalised Discounted Cumulative Gain (NDCG) [32] to quantify the quality of recommendations. These evaluations, however, do not take into account whether the observed user-item interactions were thematically driven, so they are less suitable for evaluating the quality of KG recommendations and penalise the normal behavior of our approach (i.e. not recommending thematically unrelated items). For instance, one obstacle are the many confounding factors that lead users to interact with a piece of content, such as user preference, or popularity. Therefore, we require a method that can pick up on the signal that we are most interested in, which in our case is thematic similarity.

There are two main challenges to isolating thematic signals from user-item interactions. First, it is difficult to identify the thematically relevant items-pairs, because identifying communities at the correct resolution in a multipartite, multilayer network, such as the ones we construct here, is not always practical. The second challenge is that it would be counterproductive to penalise thematically distant items that are correctly ranked lower in the recommendation list, even if they were favourably rated by the user. We assume that each item $m \in \mathcal{M}_u$ belongs to a thematic category (i.e. one of the KG’s communities), whichever that may be, and that at least some thematic categories have multiple items rated by the user. As we do not know which pairs in $z(v_m)$ are thematically relevant, we take an agnostic approach and consider only the highest ranked item favourably rated by the user. This item is likely to belong to the same thematic cluster as the seed; therefore, the position it appears in the ranked list can be used to evaluate the quality of our recommendation. In this way, we minimize the risk of penalising the algorithm for correctly down-ranking thematically distant items.

Specifically, our approach is to first select an item $m \in \mathcal{M}_u$ as a seed and generate a ranked list of recommendations $z(v_m)$, according to Eq. (4.1). The highest ranked item in $z(v_m)$ that is also relevant for our user u is:

$$(5.1) \quad \mu_u(m) = \arg \min_{m' \in \mathcal{M}_u: m' \neq m} \text{rank}(m'|z(v_m)),$$

and its rank is

$$(5.2) \quad \rho_u(m) = \min_{m' \in \mathcal{M}_u: m' \neq m} \text{rank}(m'|z(v_m)).$$

Item $\mu_u(m)$ is the candidate most likely to be thematically related to the seed m amongst the items \mathcal{M}_u that are relevant to the user u . Focusing on $\mu_u(m)$ avoids penalising the method for correctly placing thematically distant items in \mathcal{M}_u lower in the rankings. To evaluate the overall quality of recommendations for user u , we seed with every available item $m \in \mathcal{M}_u$ and compute an average score for each user which we call the *Normalised Maximum Relevance Gain (NMRG)*:

$$(5.3) \quad \text{NMRG}_u = \frac{1}{|\mathcal{M}_u|} \sum_{m \in \mathcal{M}_u} \frac{1}{\varpi_u(m)} \frac{\tau_u(\mu_u(m))}{\log_2(1 + \rho_u(m))},$$

where $\tau_u(\cdot)$ is the relevance of an item to user u (e.g. a rating or affinity score), and

$$(5.4) \quad \varpi_u(m) = \max_{m' \in \mathcal{M}_u: m' \neq m} \tau_u(m')$$

is a normalisation constant representing the maximum possible value of the score, corresponding to the case when the item with the highest relevance to the user appears in the first position of $z(v_m)$ (i.e. optimal outcome).

The NMRG considers only the position of the highest-ranked relevant item whereas the NDCG is a weighted average over the relevance of all items. Thus the NMRG is similar to the NDCG in the sense that it is proportional to the item’s relevance to the user, inversely proportional to the log-rank, and normalised by the maximum possible score. The main difference is that the NMRG is no longer cumulative, but truncated to only one item: the *highest-ranked relevant* one. Finally, note that a term in the sum of Eq. (5.3) can only be zero if $\mathcal{M}_u \cap z(v_m) = \emptyset$, when an item m has no thematically related counterparts in \mathcal{M}_u .

6. Applications. We test our thematic recommendation method in two different settings: thematic music recommendation, which we evaluate by monitoring user engagement in an AB test, and thematic movie recommendation, which we evaluate with the NMRG score. While the results of the AB test are not reproducible without access to a streaming service with a large user base, the film recommendations are performed on publicly available data, and can be reproduced by anyone.

6.1. Thematic music recommendation: AB test. We evaluate our thematic method for music recommendations. Specifically, we evaluate our KG-based recommendations in a new market for Spotify with two important characteristics that make it suitable for thematic recommendations: a) music consumption in this market is driven by thematic relationships among the entities in a KG which include people, music and movies, and b) the market comes with a new catalogue, so there is little historical user consumption data (i.e. cold market and cold catalogue). We performed an AB test where we compared the performance of playlists obtained with our KG-based

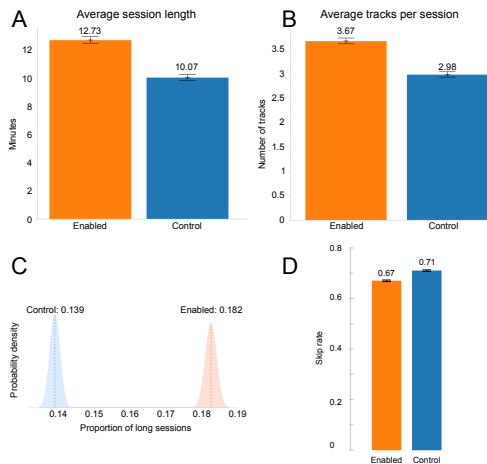


Fig. 3: AB test results. **A**: average session length in minutes. **B**: average number of tracks per session. **C**: proportion of sessions of five tracks or more. **D**: skip rate.

method (as described in Sec. 2 and 3), and the CF system that was in production at the time. Our hypothesis is that we can better capture the dynamics of consumption in this market/catalogue with our thematic approach than with the CF system. We produce song-song recommendations by creating a playlist from each seed (i.e., a track) in the following way: we first obtain 300 candidates using equation (4.1). Then, we reduce the candidate list to a 40-song playlist by re-ranking candidates to enforce artist/album separation rules, and ensure mood consistency using audio attributes. We set the value of the salience parameters by working closely with an editorial team with expertise in this market and deep knowledge of the catalogue. For the test, we assembled a cohort of 100,000 users, 50% of which received our KG-based recommendations, and a control group of 50% received the CF recommendations. After two weeks, the group exposed to KG reported (see Fig. 3):

- 26% increase in consumption time compared to control,
- 31% increase in users with sessions of 5 tracks or more,
- 6% decrease in skip rate, and
- 5% increase in save rate.

We replicated these results in two further AB tests with the same sample size, which gave us the confidence to roll out the system to all users, and this is currently the recommendation system in production in this market.

6.2. Thematic movie recommendations: offline evaluation. Next, we test our thematic KG recommendation method for movie recommendations. In this application, we combine data from two different sources: the Movie Database (TMDb) [1] metadata for the construction of a KG, and the MovieLens 20M dataset [25] as a source of explicit user feedback.

We obtained metadata from TMDb’s API [1] for all movies in the MovieLens dataset with a valid TMDb identifier. These data capture extensive information on how different types of entities (movies, people, and companies) are interrelated. For example, there are over 800 job titles available, including categories such as cast member, director, producer, and even animal wrangler. The data also include the

order in which actors appear in a movie’s credits, which we use as a proxy measure of importance of an acting credit. Additionally, the dataset contains 17k unique keywords of which 10k are associated with more than one movie. Finally, each movie in the TMDb data is associated with a popularity score, based on engagement metrics, such as votes, views, number of users who marked it as favourite, users who added it to their “watchlist”, and release date [1]. Figure 2B shows the distribution of these popularity scores, which we incorporate into the weight of the KG connections elevated to an exponent which is also a parameter. For full details on the graph construction we refer to Appendix A.1.

The MovieLens 20M dataset [25] comprises 20 million user ratings of 27k movies by 138k users; each rating has an integer value between 1 and 5. The data contain TMDb identifiers for most movies and 12 million associations of 15k movies with 1,100 tags, capturing properties of movies such as sci-fi, comedy, action, surreal, funny, classic, romance [49]. See Appendix A.2 for the details on the data preprocessing and train/test set splits.

6.2.1. Thematic signals in user data. In Sec. 5.1, we hypothesised that users display thematic preferences (such as favourite directors, actors, genres and so on); we can test this using user-item ratings in combination with our KG. If our hypothesis is correct, we expect that the path in the KG between movies rated favourably by the same user to be shorter than between two movies chosen at random. To test this hypothesis, we sample the sets of ratings (\mathcal{M}_u) of 1,000 users, and measure the shortest path in the KG between all pairs $(m, m') \in \mathcal{M}_u$. We compare these to the paths between sets of movies with the same popularity as in \mathcal{M}_u , and between movies chosen uniformly at random from the 27k movies in the dataset. Figure 2E shows the distribution of shortest paths between movies in the three sets of ratings from 1,000 users each. The distances between movies favourably rated by users are shorter than between random movies, which confirms that users do tend to interact with movies that are *thematically similar* to each other. This result shows that we can find thematic signals in user-item interaction data.

6.2.2. Parameter estimation. Our multilayer models of KGs have several free parameters including saliences and the PageRank teleportation probability, whose value may have a large impact on the quality of recommendations, and are often unknown a priori. We use the NMRG evaluation framework from Sec. 5 along with random search [5] to find parameter values that optimise the thematic relevance of recommendations, as measured by the NMRG score. The random walk model from the TMBDb data KG (Fig. 2C and D) has 14 free parameters (13 saliences and the PageRank teleportation). We sample each column of the salience matrix in Eq. (3.1) from a Dirichlet distribution with dimensions equal to the number of corresponding non-zero blocks. For example, the dimensionality of the Dirichlet random variable for the first column in Eq. (3.1) is three. We draw each column from a flat Dirichlet distribution and the teleportation parameter from the uniform $\mathcal{U}(0, 1)$ distribution. For each randomly sampled parameter configuration, we compute recommendations for every one of the 103.9k users in the training set, and average the NMRG across all users at 1, 10, and 20 recommendations. We then fine tune the teleportation probability by performing a sweep using the best configuration from the Dirichlet sampling. See Appendix B for more details on the parameter search and a discussion of our findings.

Method	NMRG@1	NMRG@10	NMRG@20
Thematic KG recs	14.83(±0.16)	30.70(±0.25)	35.07(±0.24)
Popularity baseline	0.53(±0.09)	19.77(±0.23)	22.02(±0.21)
Random seed baseline	4.91(±0.07)	16.28(±0.17)	20.86(±0.18)
Random item baseline	7.36(±0.10)	21.95(±0.21)	27.28(±0.21)
Unseeded PageRank baseline	0.70(±0.11)	23.60(±0.25)	26.95(±0.21)
TMDb keywords	4.44(±0.08)	7.34(±0.10)	8.57(±0.10)
Movielens tag genome	12.61(±0.13)	22.84(±0.18)	25.02(±0.18)
SVD	14.46(±0.14)	30.70(±0.20)	33.22(±0.19)
TMDb CF	14.36(±0.15)	28.68(±0.20)	31.14(±0.20)

Table 1: Offline evaluation NMRG scores at 1, 10, and 20 items, for all methods and baselines.

6.2.3. Comparison with baselines and alternative recommendation systems. We compare the quality of our recommendations against four baselines and three alternative recommendation methods:

- a) *Popularity baseline*: the recommendations are the same for every seed m . The ranking consists of the $K \in \{1, 10, 20\}$ most popular items in the catalogue, excluding the seed. This baseline represents the hypothesis that user engagement is 100% driven by popularity.
- b) *Random seed baseline*: we replace each seed m with a random item of *similar popularity**. This baseline represents the hypothesis that user engagement is driven by the popularity of the seed rather than the seed itself.
- c) *Random item baseline*: we retain the original seed m but replace each recommendation in $z(v_m)$ by a random item of *similar popularity*. This baseline represents the hypothesis that user engagement is driven by the popularity of the recommended items.
- d) *Unseeded baseline*: we rank items according to their unseeded PageRank (i.e. uniform teleportation). This baseline represents the hypothesis that the node centrality in the KG drives user engagement.
- e) *TMDb keywords*: we rank items based on the similarity of their associated TMDb keywords, as measured by the Dice score [6].
- f) *Movielens tag genome*: we use the graded associations between films and keywords [25] provided by the Movielens dataset to rank candidates by the Euclidean distance of their association vector and the seed’s.
- g) *SVD*: we use the open-source library **surprise** [28] to obtain a low-rank approximation of the ratings matrix using Singular Value Decomposition (SVD). In this approach, we rank items based on their dot-product similarity to the seed.
- h) *TMDb CF*: we use the collaborative filtering recommendation currently in use on the TMDb website [1].

We report NMRG scores at 1, 10, and 20 recommendations calculated on the test set only. We compare all methods using the same data, except the MovieLens tag genome and SVD, where only a subset of movies is available. We remove ratings that are not available for evaluation to avoid penalising these methods, even though

* Items of similar popularity are randomly selected from the 25 items with the closest popularity score to m .

this may lead to results slightly skewed in their favour. Table 1 contains the average NMRG score for each method with 99% confidence intervals. The performance of the baselines highlights the importance of the information encoded in the KG (in line with our findings in Fig. 2E); the best performing baselines are the ones that exploit the structure of the KG. Among the alternative methods, recommending based on TMDb keywords has the poorest performance, followed by the MovieLens tag genome. The TMDb CF recommendation performs better than all baselines and is only outperformed by the SVD approach, which is tuned on the MovieLens ratings. Our thematic recommendations performs as good as the SVD method or better in some cases. However, it is important to note that we do not claim superior performance to CF methods on *personalisation*, and the NMRG is specifically tailored to highlight the *thematic* aspects of the KG recommendations.

7. Discussion. In this work, we introduce a method for thematic item-item recommendations using knowledge graphs. We represent a KG as a multilayer network whose layers correspond to different connection types, and interlayer couplings connect nodes that represent the same entity across layers. By incorporating parameters that encode the salience of each of the connections, we are able to compare different types of connections on the same footing, allowing us to unify different connection types under the same analysis framework. Representing the KG as a multilayer network enables us to draw from a wealth of network theoretic techniques, such as random walks. We use personalised PageRank to produce recommendations that can be fine-tuned by calibrating the salience of the connections or the teleportation probability, which controls how exploratory recommendations are.

We evaluated our method for music recommendations in a specific market using an AB test. Our results show that our thematic recommendations perform significantly better than the CF in place in a range of engagement metrics. With this strong performance we rolled out the system for all users in this market. We also evaluated our framework for movie recommendations on publicly-available data. We built a KG from TMDb data and tested the performance of our method using the NMRG metric we introduce in this paper on the MovieLens 20M dataset. The evaluation framework also allows us to calibrate the parameters in the model to optimise the NMRG score and discriminate among competing models. We perform random search and parameter sweeps to identify parameter sets with good performance, and find that our method outperforms both thematic and collaborative filtering methods.

An important challenge arises during the evaluation of thematic recommendations: it is difficult to disentangle the effect of item popularity from user-item interactions. Our method performs best after incorporating the TMDb item popularity score in the connection weights. The more we accentuate the popularity, the better the model performs. However, recommending items purely based on popularity as a baseline does not perform well. This is evidence that the combination of the graph structure and popularity is required for high user satisfaction, which is also consistent with our experience in commercial applications.

Our framework is flexible and allows more complex random walk dynamics to be easily incorporated. For example, we can encode higher-order Markov chains [42, 30] at the layer level by constructing more granular roles (e.g. Tarantino acting in the 90s, Tarantino acting in the 00s and so on) and manipulating the values in the salience matrix α , although this comes at the cost of a larger system. We used the steady state of a diffusion process to recommend items. However, transients of diffusion in discrete or continuous time offer attractive alternatives. Transients do

not require irreducibility of the rate matrix, and extremely fast numerical solvers are readily available. The duration of the diffusion can be optimised globally or, more interestingly, be personalised to users, with longer durations for more adventurous users. See Ref. [37] for a recent example in discrete time.

Our framework is currently limited by the fact that it is challenging to evaluate purely thematic recommendations (i.e. without popularity in the model) due to the lack of datasets where the influence of popularity is absent. Therefore, an important task that remains is the creation of purely thematic datasets and developing evaluation methods to disentangle thematic relatedness and popularity. Finally, it is possible to integrate CF approaches into this framework by adding an item-user layer to the model of the KG. This is a promising direction of research, as it combines thematic and collaborative recommendations in a principled, transparent framework that can be understood, manipulated and optimised.

Acknowledgments. We thank Jonathan Berschadsky, Maria Dominguez, Katarzyna Drzyzga, Shahar Elisha, Martin Gould, Johnny Hunter, Nachiket Londhe, Laurence Pascall, Jyotsna Venkataramanan, and Linden Vongsathorn for their help and support as well as their useful feedback and comments. We are also thankful to the developers of TMDb for kindly providing us with support and access to their API.

Appendices

A. Movie data and preprocessing.

A.1. Construction of the knowledge graph. The TMDb Knowledge Graph includes $P = 107,143$ people, $F = 26,698$ movies, and $K = 17,592$ keywords. We construct the multilayer model of the KG with three roles for people entities: actor (**act**), director (**dir**), and producer (**prod**). We also consider a role (**desc**) for keywords describing each movie. We represent movies as a single **movie** role instead of modelling movies in four distinct passive roles (i.e. **m-dir**, **m-prod**, **m-act**, and **m-key**). As discussed in Sec. 2, merging roles reduces the size of the supra adjacency matrix (in this case by 80k rows and columns), which helps speed up computations. Figure 2D shows the sparsity profile and block structure of the supra-adjacency matrix, comprising four distinct layers. We experimented with several versions of the KG with different types of connections[†]. For brevity, we only present the configuration that performed best here.

We incorporate the popularity of the movies by letting the weight of connections to each movie m to be proportional to p_m^γ , where $\gamma \in \mathbb{R}$ is an exponent controlling the importance of popularity. For instance, $\gamma > 1$ accentuates differences in popularity, $\gamma = 0$ discards them, and $\gamma < 0$ reverses their effect. We set the weight of **directs**, **produces** and **describes** connections to p_m^γ . The weights of **acts_in** connections are set to $\frac{p_m^\gamma}{\log_2(c_i+1)}$, where c_i is the position of actor i in the credits, to account for both popularity and credit order. We observed that higher values of the popularity exponent γ result in increased performance, and fixed this parameter to a large value ($\gamma = 30$). This highlights the overwhelming importance of popularity, and in effect sparsifies the network by connecting each node to its' most popular associated movie only.

[†]Including animal wrangler, of course.

Step	Description	#ratings	#movies
#1	Remove ratings without TMDb id	19,987,681	26,483
#2	Remove duplicates (keeping latest)	19,987,649	26,483
#3	Remove movies without metadata	19,950,334	26,185
#4	Remove ratings below user median	13,032,003	23,431
#5	Remove movies with 2 or fewer ratings	13,028,453	19,881
#6	Keep only top 250 ratings per user	10,563,717	18,319

Table 2: Filtering steps applied to the dataset, with their corresponding impact on the number of ratings and number of movies.

	Actor	Movie	Keyword	Producer	Director
Actor	-	0.798	-	0.643	0.149
Movie	0.3585	-	1	0.057	0.424
Keyword	-	0.012	-	-	-
Producer	0.550	0.064	-	-	0.427
Director	0.091	0.126	-	0.299	-

Table 3: Best salience values discovered using random search.

A.2. MovieLens 20M dataset. To obtain the data we use in our experiments, we apply heuristic filtering steps to the user ratings in the order shown in Table 2. The total number of users is not affected by the filtering. More specifically, we retain only favourable user ratings (above the user’s median rating) and remove the rest to avoid rewarding methods that recommend items the users may not have liked. We also remove movies with 2 ratings or fewer, because we expect those ratings to be noisy. Finally, we keep only the 250 maximum ratings per user, so that we do not bias our metrics in favour of power users.

We split the user set uniformly at random into a train (75%) and test set (25%). We use the train set for random parameter search and training. We report results on the test set only for all methods. The resulting train set comprises 103.9k users, 17.6k movies and 7.9M ratings, while the test set consists of 34.6k users, 15.1k movies and 2.6M ratings.

B. Random parameter search. We use the evaluation framework introduced in Sec. 5 along with random search [5] to find good values of the 14 free parameters in our model. The teleportation parameter ρ in Eq. 3.4, takes positive values between $(0, 1)$, so we sample it from uniform distribution $\rho \sim \mathcal{U}(0, 1)$. We sample the values of each column in Eq. 3.1 from a Dirichlet distribution with unit concentration parameter $\alpha_i \sim \text{Dir}(\mathbf{1})$, where i is the column index.

Evaluating on the full train set requires a considerable amount of computation, despite the fact that the process can easily be performed in parallel. To optimise computations, we use a two-stage approach: we first compute ranked lists $z(v_m)$ by seeding at each one of the 17.6k movies in the training set, then we use these lists to evaluate the NMRG for each user following Eq. (5.3). Such pre-computing speeds up our calculations significantly. We present the best salience parameters as identified by the random search in Table 3.

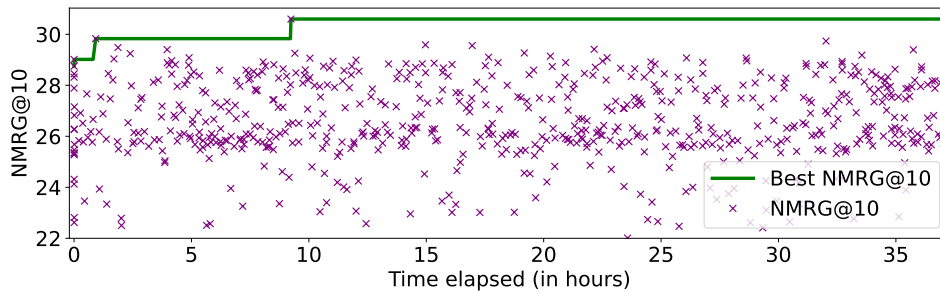


Fig. 4: Random search NMRG@10 and best NMRG@10 as the experiment evolves.

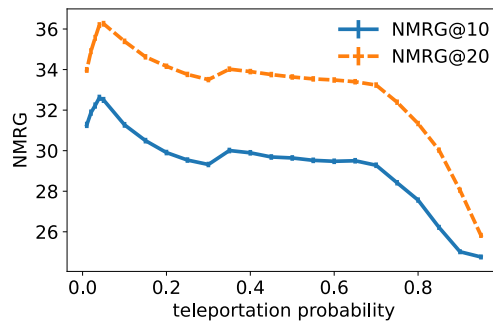


Fig. 5: Teleportation parameter (ρ) sweep.

The results of our random parameter search were obtained after evaluating 655 random parameter configurations. For this computation we employed 5 machines with 96 vCPUs each; the total run time was approximately 37 hours. In Fig. 4 we present the NMRG@10 and best NMRG@10 as the experiment evolves.

Figure 5 shows a sweep of the teleportation probability ρ . The optimal value of the NMRG is achieved when $\rho \approx 0.12$.

REFERENCES

- [1] *The Movie Database (TMDb)*, 2008, <https://www.themoviedb.org/> (accessed 2020-05-27).
- [2] Q. AI, V. AZIZI, X. CHEN, AND Y. ZHANG, *Learning heterogeneous knowledge base embeddings for explainable recommendation*, *Algorithms*, 11 (2018), p. 137.
- [3] M. BEGUERISSE-DÍAZ, M. A. PORTER, AND J.-P. ONNELA, *Competition for popularity in bipartite networks*, *Chaos*, 20 (2010), p. 043101, <https://doi.org/10.1063/1.3475411>.
- [4] L. BELLOMARINI, E. SALLINGER, AND G. GOTTLÖB, *The vadalog system: Datalog-based reasoning for knowledge graphs*, *Proc. VLDB Endow.*, 11 (2018), pp. 975–987, <https://doi.org/10.14778/3213880.3213888>.
- [5] J. BERGSTRA AND Y. BENGIO, *Random search for hyper-parameter optimization*, *JMLR*, 13 (2012), pp. 281–305.
- [6] A. CARASS, S. ROY, A. GHERMAN, J. C. REINHOLD, A. JESSON, T. ARBEL, O. MAIER, H. HANDELS, M. GHAFOORIAN, B. PLATEL, ET AL., *evaluating white matter lesion segmentations with refined sorensen-dice analysis*, *Scientific Reports*, 10 (2020), pp. 1–19.
- [7] Y.-S. CHANG, F. NIE, Z. LI, X. CHANG, AND H. HUANG, *Refined spectral clustering via embedded label propagation*, *Neural Computation*, 29 (2017), pp. 3381–3396, <https://doi.org/10.1162/NEUR-2016-0183>.

- 1162/neco.a.01022.
- [8] S. CHAUDHARI, A. AZARIA, AND T. MITCHELL, *An entity graph based recommender system*, AI Communications, 30 (2017), pp. 141–149.
 - [9] X. CHEN, S. JIA, AND Y. XIANG, *A review: Knowledge reasoning over knowledge graph*, Expert Systems with Applications, 141 (2020), p. 112948, <https://doi.org/10.1016/j.eswa.2019.112948>.
 - [10] P. COVINGTON, J. ADAMS, AND E. SARGIN, *Deep neural networks for youtube recommendations*, in Rec. Sys., 2016, <https://doi.org/10.1145/2959100.2959190>.
 - [11] H. CRAMER, J. GARCIA-GATHRIGHT, A. SPRINGER, AND S. REDDY, *Assessing and addressing algorithmic bias in practice*, Interactions, 25 (2018), pp. 58–63.
 - [12] M. DE DOMENICO, A. SOLÉ-RIBALTA, S. GÓMEZ, AND A. ARENAS, *Navigability of interconnected networks under random failures*, PNAS, 111 (2014), pp. 8351–8356, <https://doi.org/10.1073/pnas.1318469111>.
 - [13] G. DEL CORSO, A. GULLI, AND F. ROMANI, *Fast pagerank computation via a sparse linear system*, Internet Mathematics, 2 (2005), pp. 251–273.
 - [14] R. DIESTEL, *Graph Theory*, Springer, 2005.
 - [15] X. DONG, E. GABRILOVICH, G. HEITZ, W. HORN, N. LAO, K. MURPHY, T. STROHMANN, S. SUN, AND W. ZHANG, *Knowledge vault: A web-scale approach to probabilistic knowledge fusion*, in KDD, 2014, pp. 601–610.
 - [16] X. L. DONG, *Challenges and innovations in building a product knowledge graph*, in KDD, 2018, pp. 2869–2869.
 - [17] Y. DONG, N. V. CHAWLA, AND A. SWAMI, *Metapath2vec: Scalable representation learning for heterogeneous networks*, in KDD, 2017, pp. 135–144, <https://doi.org/10.1145/3097983.3098036>.
 - [18] R. EBERT, *Amores perros*, April 2001, <https://www.rogerebert.com/reviews/amores-perros-2001>.
 - [19] L. EHRLINGER AND W. WÖSS, *Towards a definition of knowledge graphs.*, SEMANTICS, 48 (2016).
 - [20] C. EKSOMBATCHAI, P. JINDAL, J. Z. LIU, Y. LIU, R. SHARMA, C. SUGNET, M. ULRICH, AND J. LESKOVEC, *Pixie: A system for recommending 3+ billion items to 200+ million users in real-time*, in WWW, 2018, pp. 1775–1784.
 - [21] A. GILOTTE, C. CALAUZÈNES, T. NEDELEC, A. ABRAHAM, AND S. DOLLÉ, *Offline a/b testing for recommender systems*, in WSDM, 2018, pp. 198–206.
 - [22] D. F. GLEICH, *PageRank beyond the Web*, SIAM Review, 57 (2015), pp. 321–363.
 - [23] J. GOPE AND S. K. JAIN, *A survey on solving cold start problem in recommender systems*, in Int. Conf. on Computing, Communication & Automation, 2017, pp. 133–138.
 - [24] A. GRUSON, P. CHANDAR, C. CHARBUILLET, J. MCINERNEY, S. HANSEN, D. TARDIEU, AND B. CARTERETTE, *Offline evaluation to make decisions about playlist recommendation algorithms*, in WSDM, 2019, pp. 420–428.
 - [25] F. HARPER AND J. KONSTAN, *The movielens datasets: History and context*, AcM transactions on interactive intelligent systems (tiis), 5 (2015), pp. 1–19.
 - [26] D. HIGHAM AND A. TAYLOR, *The sleekest link algorithm*, Mathematics Today, 39 (2003), pp. 192–197.
 - [27] A. HOGAN, E. BLOMQUIST, M. COCHEZ, C. D’AMATO, G. DE MELO, C. GUTIERREZ, J. E. L. GAYO, S. KIRrane, S. NEUMAIER, A. POLLERES, R. NAVIGLI, A.-C. N. NGOMO, S. RASHID, A. RULA, L. SCHMELZEISEN, J. SEQUEDA, S. STAAB, AND A. ZIMMERMANN, *Knowledge graphs*, arXiv:2003.02320, (2020), <https://arxiv.org/abs/2003.02320>.
 - [28] N. HUG, *Surprise, a Python library for recommender systems*. <http://surpriselib.com>, 2017.
 - [29] M. KIVELÄ, A. ARENAS, M. BARTHELEMY, J. P. GLEESON, Y. MORENO, AND M. A. PORTER, *Multilayer networks*, Journal of Complex Networks, 2 (2014), pp. 203–271, <https://doi.org/10.1093/comnet/cnu016>.
 - [30] R. LAMBIOTTE, M. ROSVALL, AND I. SCHOLTES, *From networks to optimal higher-order models of complex systems*, Nature physics, 15 (2019), pp. 313–320.
 - [31] L. LÜ, M. MEDO, C. H. YEUNG, Y.-C. ZHANG, Z.-K. ZHANG, AND T. ZHOU, *Recommender systems*, Physics Reports, 519 (2012), pp. 1 – 49, <https://doi.org/https://doi.org/10.1016/j.physrep.2012.02.006>. Recommender Systems.
 - [32] C. D. MANNING, P. RAGHAVAN, AND H. SCHÜTZE, *Introduction to Information Retrieval*, CUP, 2008, <https://doi.org/10.1017/CBO9780511809071>.
 - [33] T. MARTIN, X. ZHANG, AND M. E. J. NEWMAN, *Localization and centrality in networks*, Phys. Rev. E, 90 (2014), p. 052808, <https://doi.org/10.1103/PhysRevE.90.052808>.
 - [34] N. MASUDA, M. PORTER, AND R. LAMBIOTTE, *Random walks and diffusion on networks*, Physics Reports, 716–717 (2017), pp. 1–58, <https://doi.org/10.1016/j.physrep.2017.07.007>.

- [35] J. MCINERNEY, B. LACKER, S. HANSEN, K. HIGLEY, H. BOUCHARD, A. GRUSON, AND R. MEHROTRA, *Explore, exploit, and explain: personalizing explainable recommendations with bandits*, in Rec. Sys., 2018, pp. 31–39.
- [36] F. MONTI, M. BRONSTEIN, AND X. BRESSON, *Geometric matrix completion with recurrent multi-graph neural networks*, in NIPS, 2017, pp. 3697–3707.
- [37] A. N. NIKOLAKOPOULOS, D. BERBERIDIS, G. KARYPIS, AND G. B. GIANNAKIS, *Personalized diffusions for top-n recommendation*, in Rec. Sys., 2019, pp. 260–268, <https://doi.org/10.1145/3298689.3346985>.
- [38] L. PAGE, S. BRIN, R. MOTWANI, AND T. WINOGRAD, *The PageRank citation ranking: Bringing order to the Web*, in WWW, 1998, pp. 161–172.
- [39] F. POURSAZBI-SANGDEH, D. G. GOLDSTEIN, J. M. HOFMAN, J. W. VAUGHAN, AND H. WALLACH, *Manipulating and measuring model interpretability*, arXiv, 1802.07810 (2018).
- [40] F. RICCI, L. ROKACH, B. SHAPIRA, AND P. B. KANTOR, *Recommender Systems Handbook*, Springer, 2010.
- [41] A. ROSSI, D. FIRMANI, A. MATINATA, P. MERIALDO, AND D. BARBOSA, *Knowledge graph embedding for link prediction: A comparative analysis*, arXiv, 2002.00819 (2020).
- [42] M. ROSVALL, A. ESQUIVEL, A. LANCICHINETTI, J. WEST, AND R. LAMBIOTTE, *Memory in network flows and its effects on spreading dynamics and community detection*, Nature Communications, 5 (2014), pp. 1–13.
- [43] G. SALHA, R. HENNEQUIN, AND M. VAZIRGIANNIS, *Keep it simple: Graph autoencoders without graph convolutional networks*, arXiv:1910.00942, (2019), <https://arxiv.org/abs/1910.00942>.
- [44] Y. SHI, M. LARSON, AND A. HANJALIC, *Collaborative filtering beyond the user-item matrix: A survey of the state of the art and future challenges*, Comput. Surv., 47 (2014), <https://doi.org/10.1145/2556270>.
- [45] D. SIROKER AND P. KOOMEN, *A/B Testing: The Most Powerful Way to Turn Clicks Into Customers*, Wiley, 2013.
- [46] Z. SUN, J. YANG, J. ZHANG, A. BOZZON, L.-K. HUANG, AND C. XU, *Recurrent knowledge graph embedding for effective recommendation*, in Rec. Sys., New York, NY, USA, 2018, pp. 297–305, <https://doi.org/10.1145/3240323.3240361>.
- [47] M. SZOMSZOR, C. CATTUTO, H. ALANI, K. O’HARA, A. BALDASSARRI, V. LORETO, AND V. D. SERVEDIO, *Folksonomies, the semantic web, and movie recommendation*, in European Semantic Web Conf., June 2007.
- [48] D. TAYLOR, *Multiplex markov chains*, arXiv, 2004.12820 (2020), <https://arxiv.org/abs/2004.12820>.
- [49] J. VIG, S. SEN, AND J. RIEDL, *The tag genome: Encoding community knowledge to support novel interaction*, ACM Trans. Interact. Intell. Syst., 2 (2012), <https://doi.org/10.1145/2362394.2362395>.
- [50] Z. WANG, Y. TAN, AND M. ZHANG, *Graph-based recommendation on social networks*, in Int. Asia-Pacific Web Conf., 2010, pp. 116–122, <https://doi.org/10.1109/APWeb.2010.60>.
- [51] Y. XIAN, Z. FU, S. MUTHUKRISHNAN, G. DE MELO, AND Y. ZHANG, *Reinforcement knowledge graph reasoning for explainable recommendation*, in SIGIR, 2019, pp. 285–294.
- [52] S. YANG, F. HAN, Y. WU, AND X. YAN, *Fast top-k search in knowledge graphs*, in 32nd Int. Conf. on Data Engineering, 2016, pp. 990–1001.
- [53] R. YING, R. HE, K. CHEN, P. EKSOMBATCHAI, W. HAMILTON, AND J. LESKOVEC, *Graph convolutional neural networks for web-scale recommender systems*, in KDD, 2018, pp. 974–983.
- [54] X. YU, X. REN, Y. SUN, Q. GU, B. STURT, U. KHANDELWAL, B. NORICK, AND J. HAN, *Personalized entity recommendation: A heterogeneous information network approach*, in WSDM, 2014, pp. 283–292, <https://doi.org/10.1145/2556195.2556259>.
- [55] G. ZHENG, F. ZHANG, Z. ZHENG, Y. XIANG, N. J. YUAN, X. XIE, AND Z. LI, *DRN: A deep reinforcement learning framework for news recommendation*, in WWW, 2018, pp. 167–176, <https://doi.org/10.1145/3178876.3185994>.