

# Pixie-Inspired Random Walk Recommendation Algorithms

## 1. Introduction

Pixie-inspired recommendation algorithms belong to a broader family of graph-based recommenders that operate directly on a user–item interaction graph. Instead of only relying on explicit similarity calculations between users or items, these methods treat historical interactions (such as ratings, clicks, pins, purchases, or views) as edges in a bipartite graph. One partition of the graph represents users, and the other partition represents items. A connection between a user and an item indicates that the user interacted with that item in some way.

The core idea behind Pixie is to generate recommendations by launching many short, biased random walks from one or more seed nodes on this graph. These seeds can be a user profile, a recently viewed item, or a collection of items representing the user’s current session. As the random walks traverse the graph, they repeatedly land on item nodes. Items that are visited more frequently across these walks are considered stronger candidates for recommendation. In this way, the algorithm naturally exploits both direct and higher-order connections in the graph without explicitly computing full-pair similarity matrices.

In the context of this assignment, we implement a simplified, educational version of a Pixie-inspired algorithm for the MovieLens 100K dataset. We construct a bipartite graph whose nodes are users and movies, and whose edges correspond to observed ratings. From a starting user node, we perform a random walk over the graph, count how often movie nodes are visited, and return the most frequently visited movies that the user has not already rated. Although this implementation does not incorporate all of the industrial optimizations of production Pixie systems, it captures the essential conceptual idea.

## 2. Graph Construction and Normalization

The first step in a Pixie-style system is to construct an appropriate graph representation. In our case, we take the ratings DataFrame, which contains user IDs, movie IDs, and rating values, and merge it with the movies DataFrame to attach human-readable titles. We then aggregate ratings so that each (user, movie) pair has a single numeric rating, typically the mean if multiple ratings exist. To reduce user-specific rating bias (some users consistently rate higher or lower than others), we normalize ratings per user by subtracting each user’s mean rating from their individual ratings. This centers each user’s rating distribution around zero.

With normalized data, we build a bipartite adjacency list representation. We create a Python dictionary called graph in which keys are node identifiers (user IDs and movie IDs), and values are sets of neighboring node IDs. For each row in the ratings DataFrame, we add undirected edges: the user is connected to the movie, and the movie is connected back to the user. This representation is memory-efficient, simple to traverse, and avoids the need for external graph libraries. It also reflects the core structure exploited by Pixie-style methods:

connectivity between users and items through interaction edges.

While our implementation uses a basic undirected bipartite graph, industrial systems often enhance the graph with edge weights (for example, based on rating magnitude, recency, or interaction type) and additional node types (such as topics, categories, or boards). These enhancements allow more nuanced random-walk behavior and better personalization, but the fundamental concept remains the same: recommendations emerge from the structure and dynamics of the graph.

### 3. Random Walk Dynamics

Once the graph is built, the Pixie-inspired algorithm performs random walks to discover relevant items for a target user. The walk begins at the user node corresponding to the current user. At each step, we examine the neighbors of the current node, randomly select one neighbor, and move to that node. If the current node is a user, its neighbors are movies the user has rated; if it is a movie, its neighbors are users who have rated that movie. This alternating traversal between users and movies allows the walk to explore chains of “user → movie → user → movie” relationships.

During the walk, we focus on the times when we land on movie nodes. Each time a movie node is visited, we increment a counter associated with that movie. After a fixed number of steps (the `walk_length` parameter), we stop the walk and inspect the accumulated visit counts. Movies with higher visit counts are interpreted as being more relevant or more strongly connected to the starting user in the graph. Finally, we filter out movies the user has already rated and select the top-N movies by visit count as recommendations.

In our assignment implementation, the neighbor selection at each step is uniform: every neighbor of the current node has equal probability of being chosen. A more “weighted” Pixie algorithm would instead draw neighbors with probabilities proportional to some weight, such as the absolute value of the normalized rating, the recency of interaction, or a learned importance score. This weighting biases the random walk toward stronger or more relevant connections, improving the quality of recommendations.

### 4. Why Random Walks Work for Recommendations

Random walks are powerful in recommendation settings because they naturally capture complex, higher-order relationships in the data. Traditional collaborative filtering methods often rely on direct co-occurrence or similarity between users or items. For example, user-based collaborative filtering looks for users with similar rating vectors, and item-based collaborative filtering looks for items with similar rating patterns. While effective, these methods do not always capture indirect relationships, especially in sparse datasets where many users rate only a small subset of items.

A random walk on a user–item graph, by contrast, can traverse multiple hops away from the starting user or item. Consider two movies that are never rated by exactly the same users but are connected by a chain of intermediate users and items. A random walk can explore such paths, allowing the algorithm to recommend items that are indirectly related but still relevant. In other words, the random walk encodes a notion of graph proximity that is richer than simple

similarity measures.

Moreover, random walks can be tuned via parameters such as walk length, restart probability, and edge weights to balance exploration and exploitation. Short walks tend to stay close to the starting node and yield highly personalized recommendations, while longer walks explore more of the graph and may surface more diverse or serendipitous items. Production systems such as Pinterest's Pixie use many short, biased walks with various seeds and advanced pruning strategies to achieve both scalability and recommendation quality at web scale.

## 5. Real-World Applications and Extensions

Pixie-inspired algorithms have been successfully deployed in real-world systems, particularly in large-scale content discovery platforms. Pinterest popularized Pixie for recommending pins, boards, and users by modeling interactions as a rich multi-type graph and running optimized random walks that can respond to user requests in real time. Similar ideas appear in systems for recommending products, news articles, songs, and videos, where user-item graphs are large, dynamic, and highly sparse.

Beyond content recommendation, random-walk-based approaches are used in link prediction, anomaly detection, graph-based search ranking, and social network analysis. In each of these domains, the central idea is similar: use the pattern of connections in a graph to infer which nodes are most relevant or important for a particular query or starting point.

Our assignment implementation is intentionally simplified: we use a single bipartite graph with users and movies, a single unweighted random walk per query, and a modest walk length. Nevertheless, it demonstrates the core strengths of Pixie-inspired methods: the ability to leverage graph structure, discover higher-order relationships, and generate recommendations without explicitly computing dense similarity matrices. Future extensions could introduce weighted edges, multiple walks, restarts back to the starting user, or hybridization with collaborative filtering scores to further improve accuracy and robustness.