# Special Topics in Applied Mathematics I Solution 2

## XIA JIAHAN

### February 5, 2026

## 1 Introduction to Recommendation Systems

In the information explosion era, recommendation systems (RS) are essential for filtering vast data and delivering personalized content across platforms like Amazon, Netflix, YouTube, and social media, enhancing user experience and business value by predicting preferences. They use machine learning to model similarities between items and user interests, powering two main types of suggestions: personalized homepage feeds (unique to each user) and related item recommendations (e.g., suggesting science apps when viewing a math app). This enables platforms like YouTube and the Google Play Store to anticipate what to show next. The core motivation is to help users navigate immense content libraries, where search alone falls short, by surfacing relevant, sometimes unexpected items and facilitating discovery.

A recommender system matches a user's query (context) to recommendable entities (items) by learning embeddings that place queries and items in a shared vector space for efficient similarity computation. It then performs candidate generation to quickly narrow a huge corpus to hundreds or thousands of candidates, scoring to rank this smaller set and select roughly ten items to display, and finally re-ranking to apply constraints such as removing content the user dislikes, boosting fresher items, and ensuring diversity and fairness. We will discuss each stage with examples from systems like YouTube.



Figure 1: Overview of the Recommender System Process

## 2 Candidate generation

### 2.1 Overview

Candidate generation is the first stage of recommendation. Given a query, the system generates a set of relevant candidates. The following table shows two common candidate generation approaches:

| Type | Definition | Example |
|------|-----------|---------|
| content-based filtering | Uses *similarity between items* to recommend items similar to what the user likes. | If user A watches two cute cat videos, then the system can recommend cute animal videos to that user. |
| collaborative filtering | Uses *similarities between queries and items simultaneously* to provide recommendations. | If user A is similar to user B, and user B likes video 1, then the system can recommend video 1 to user A (even if user A hasn't seen any videos similar to video 1). |

Both content-based and collaborative filtering embed items and queries (contexts) into a shared low-dimensional space $E = \mathbb{R}^d$. Candidate generation then retrieves items whose embeddings are most similar to the query embedding: a similarity function $s : E \times E \to \mathbb{R}$ ranks candidates by $s(q, x)$. Common choices include cosine similarity, dot product, and Euclidean distance.

Unlike cosine similarity, dot product is sensitive to embedding norms, so large-norm vectors can be preferred even when the angle is similar. This can affect recommendations as follows:

- Popular items often get large norms, so dot product can capture popularity but may cause popular items to dominate. A tempered variant is $s(q, x) = \|q\|^\alpha \|x\|^\alpha \cos(q, x)$ for $\alpha \in (0, 1)$.

- Rare items may be updated infrequently; if initialized with large norms they can be over-recommended. To avoid this problem, be careful about embedding initialization and use appropriate regularization.