

Building Recommendation Systems in Spark ML



Janani Ravi

CO-FOUNDER, LOONYCORN

www.loonycorn.com

Overview

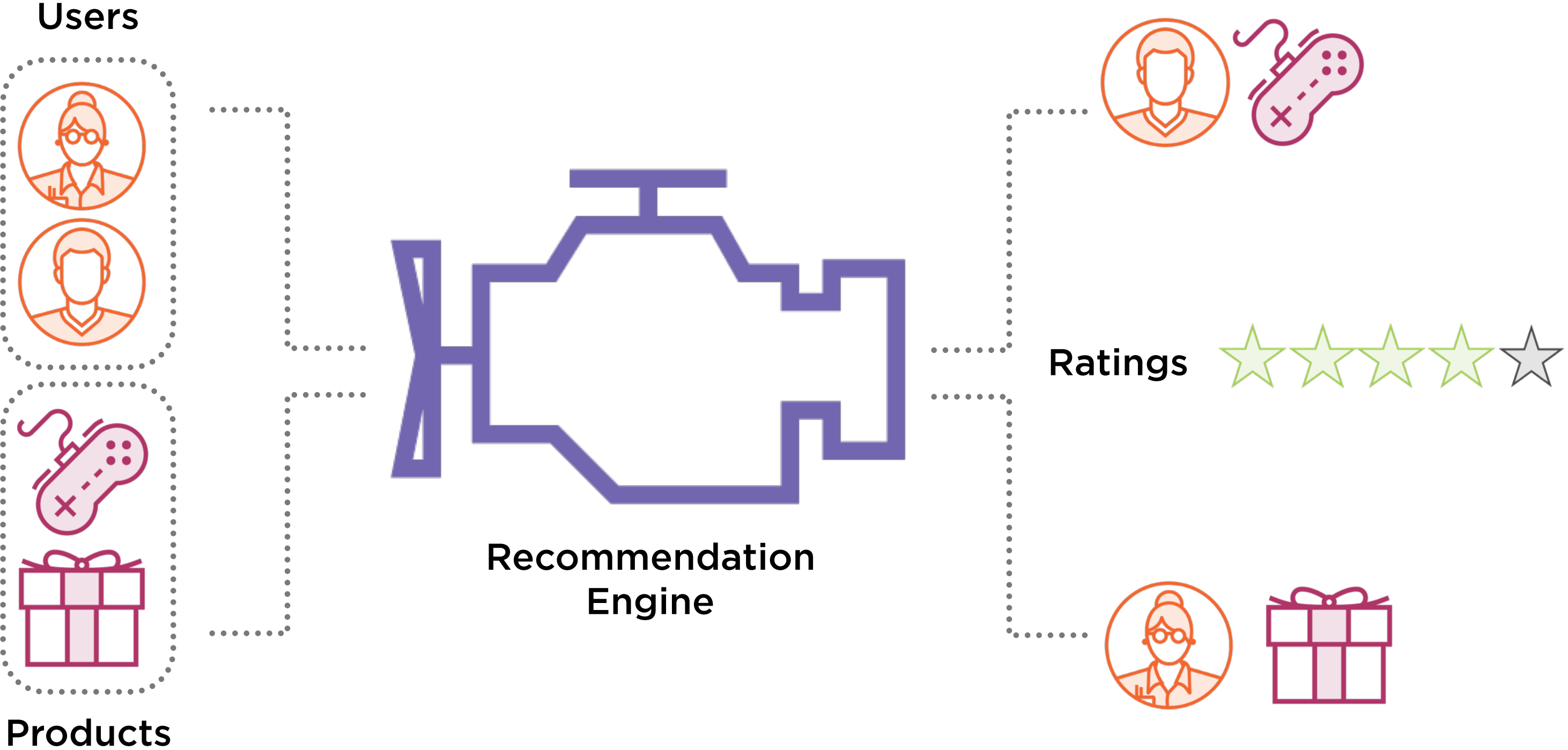
Collaborative filtering algorithms are used to make product recommendations to users

Spark offers estimators which use the ALS method to implement collaborative filtering

Abstracts you from the math involved in ALS

Recommendations can be based on explicit and implicit ratings

Recommendation Systems



Approaches to Recommendations

Content-based

Estimate rating using this user and this product alone

Collaborative

Employ information about other users, products too

Hybrid

Combine both content-based and collaborative filtering

Approaches to Recommendations

Content-based

Estimate rating using this user and this product alone

Collaborative

Employ information about other users, products too

Hybrid

Combine both content-based and collaborative filtering

Content-based Filtering

**Individual
Users**



Products



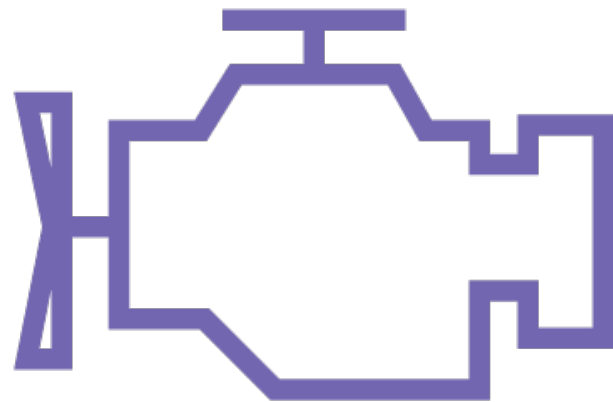
**Personalized
Recommendations**



Views

Purchases





Content-based Filtering

Match product description to user profile

Two significant drawbacks

- **Requires accurate, rich product metadata**
- **Hard to extend across product types**

Approaches to Recommendations

Content-based

Estimate rating using this user and this product alone

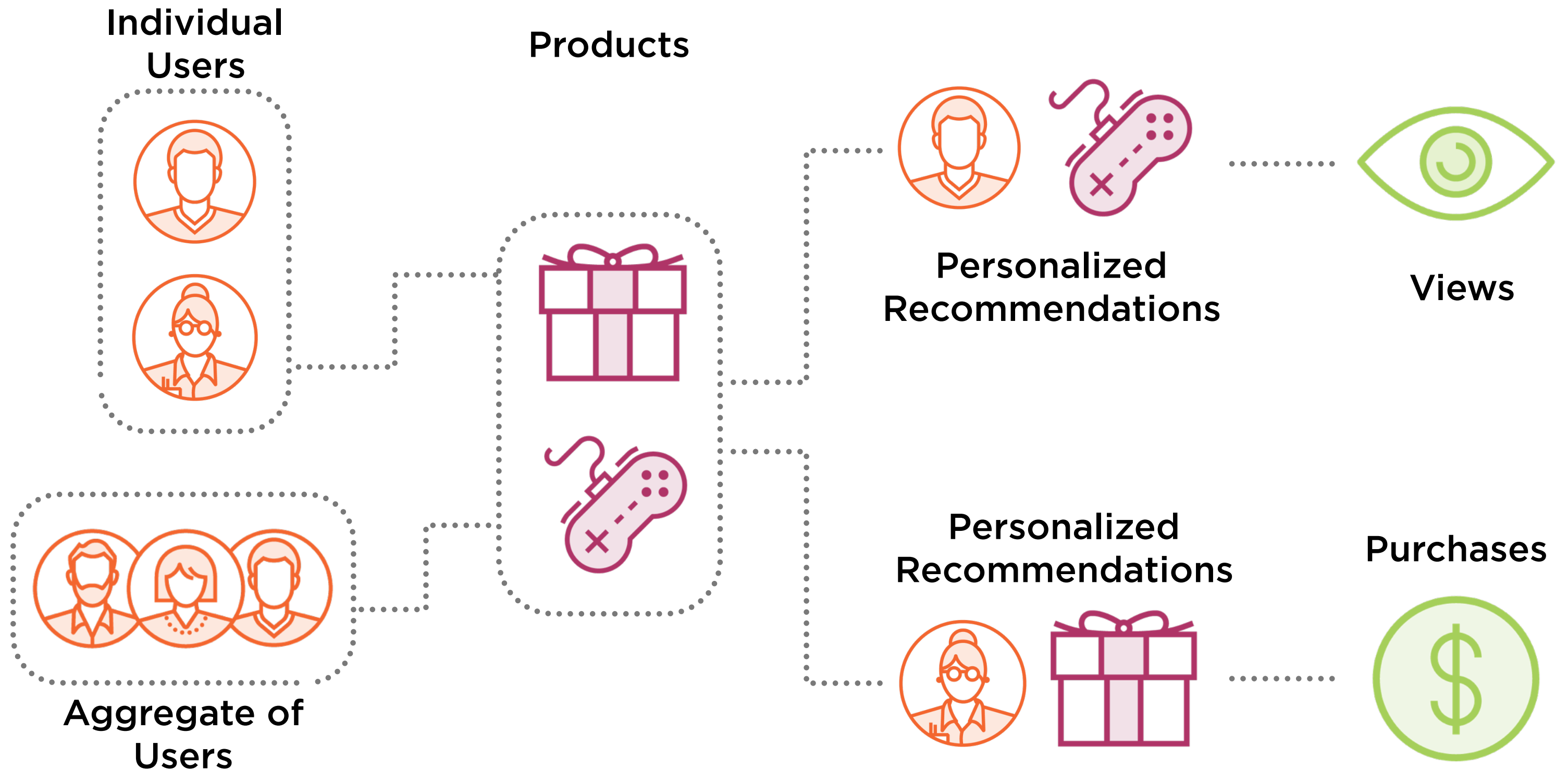
Collaborative

Employ information about other users, products too

Hybrid

Combine both content-based and collaborative filtering

Collaborative Filtering



Collaborative Filtering



Collaborative Filtering



Collaborative Filtering

Users who agreed in the past will agree in the future,
and that they will like similar kinds of items as they liked
in the past

Collaborative Filtering

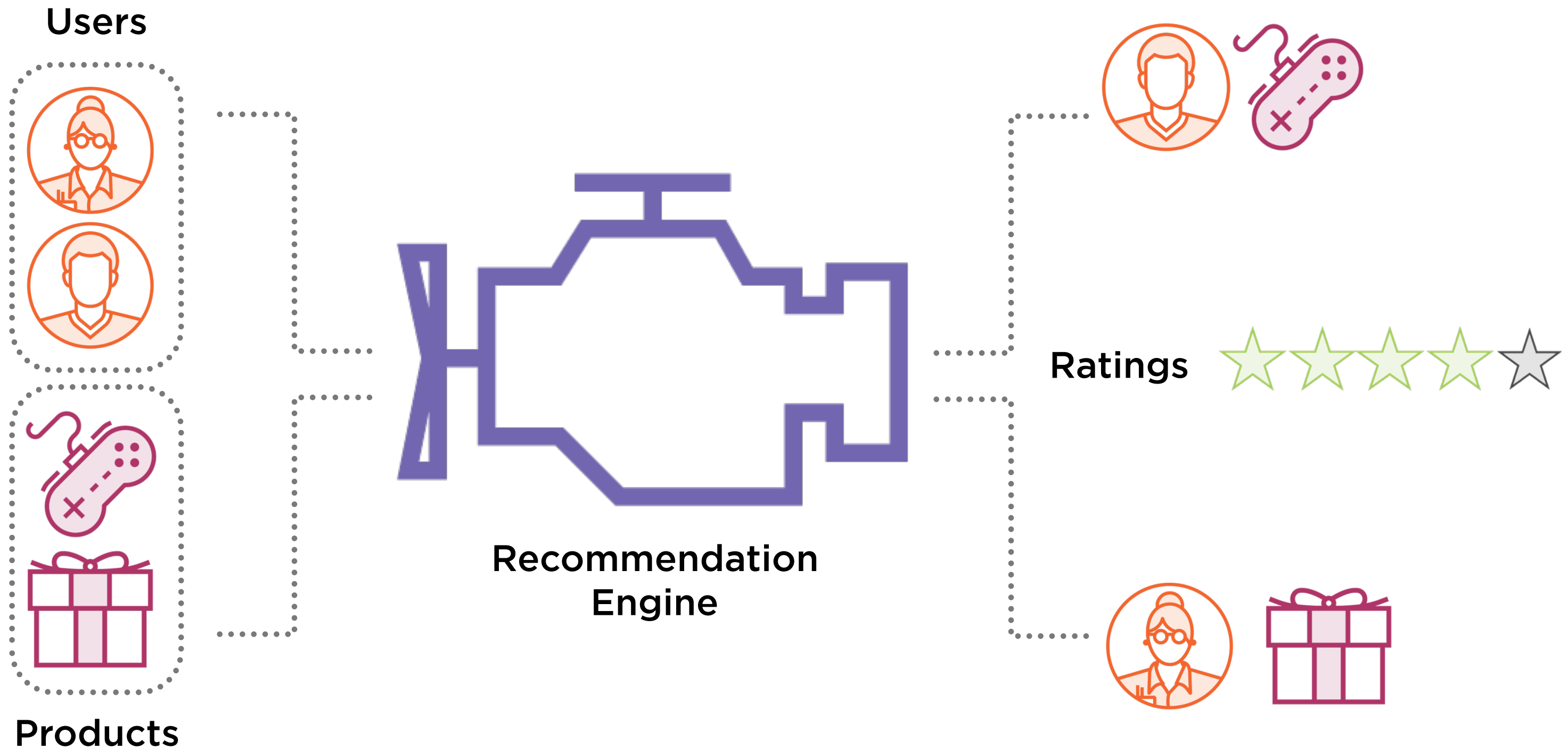
Users who agreed in the past will agree in the future,
and that they will like similar kinds of items as they liked
in the past

Collaborative Filtering

Users who agreed in the past will agree in the future,
and that they will like similar kinds of items as they liked
in the past

“People who buy X also buy Y”

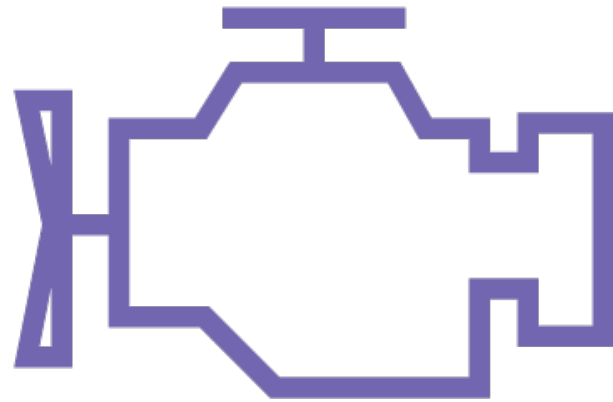
Recommendation Systems



Estimate how a user would
rate every product

Recommend the products to the
user which have the highest
estimated ratings

Ratings



Desired output of Recommendation Engine

- **Ratings Matrix:** score for each combination of user and product
- **Number of rows** = Number of users (n_u)
- **Number of columns** = Number of products (n_p)

Ratings Matrix



r_{11}	r_{12}	r_{13}	\dots	r_{1n_p}
r_{21}	r_{22}	r_{23}	\dots	r_{2n_p}
r_{31}	r_{32}	r_{33}	\dots	r_{3n_p}
\dots	\dots	\dots	r_{ij}	\dots
r_{n_u1}	r_{n_u2}	r_{n_u3}	\dots	$r_{n_un_p}$

n_p columns

n_u rows

Each element predicts how much a particular user will like a particular product

Ratings Matrix

	Product 1	Product 2	Product 3		Product n_p
User 1 	r_{11}	r_{12}	r_{13}		r_{1n_p}
	r_{21}	r_{22}	r_{23}		r_{2n_p}
	r_{31}	r_{32}	r_{33}	...	r_{3n_p}
	r_{ij}	...
	r_{n_u1}	r_{n_u2}	r_{n_u3}		$r_{n_un_p}$

Each row represents the preferences of 1 user for different products

Ratings Matrix

	Product 1	Product 2	Product 3		Product n_p
	r_{11}	r_{12}	r_{13}		r_{1n_p}
User 2 	r_{21}	r_{22}	r_{23}		r_{2n_p}
	r_{31}	r_{32}	r_{33}	...	r_{3n_p}
	r_{ij}	...
	r_{n_u1}	r_{n_u2}	r_{n_u3}		$r_{n_un_p}$

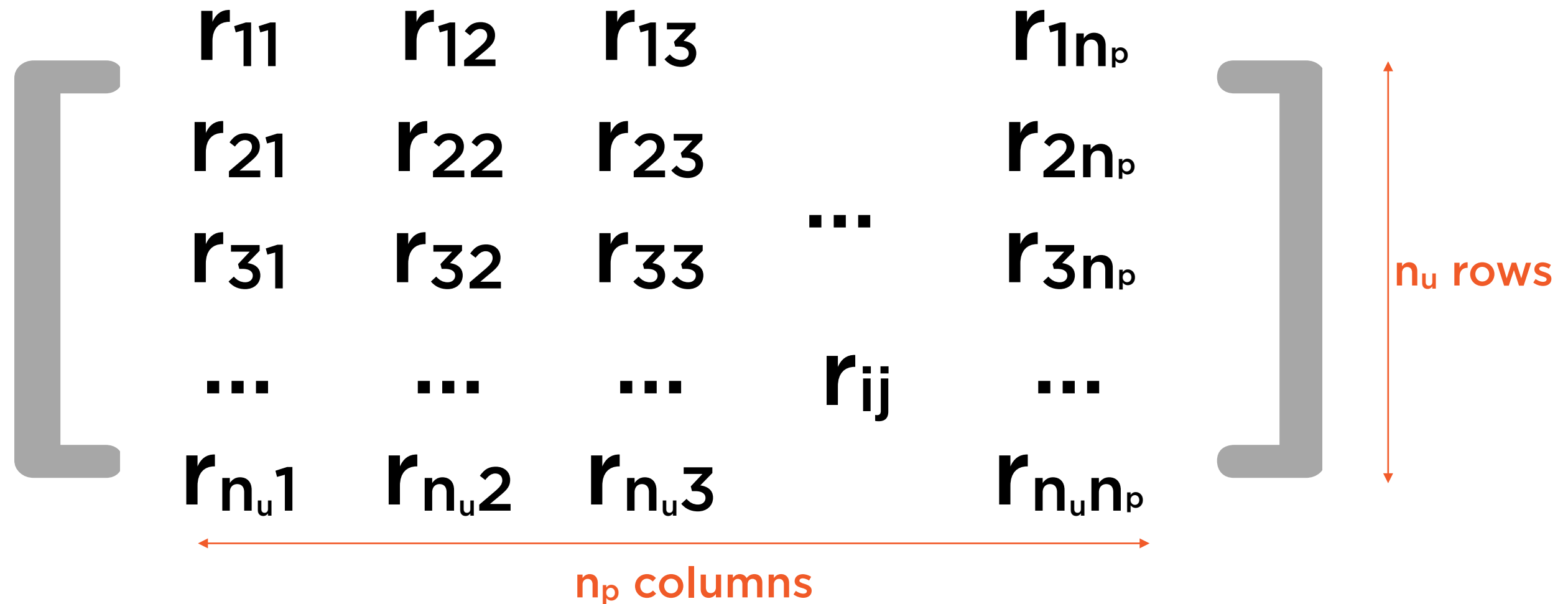
Each row represents the preferences of 1 user for different products

Ratings Matrix

	Product 1	Product 2	Product 3		Product n_p
	r_{11}	r_{12}	r_{13}		r_{1n_p}
	r_{21}	r_{22}	r_{23}		r_{2n_p}
	r_{31}	r_{32}	r_{33}	...	r_{3n_p}
	r_{ij}	...
User n_u 	r_{n_u1}	r_{n_u2}	r_{n_u3}		$r_{n_un_p}$

Each row represents the preferences of 1 user for different products

Ratings Matrix



Each column represents the preference for a single product across all users

Ratings Matrix

	Product 1				
User 1	r_{11}	r_{12}	r_{13}		r_{1n_p}
User 2	r_{21}	r_{22}	r_{23}		r_{2n_p}
	r_{31}	r_{32}	r_{33}	...	r_{3n_p}
	r_{ij}	...
User n_u	r_{n_u1}	r_{n_u2}	r_{n_u3}		$r_{n_un_p}$



Each column represents the preference for a single product across all users

Ratings Matrix

		Product 2			
User 1	r_{11}	r_{12}	r_{13}		r_{1n_p}
User 2	r_{21}	r_{22}	r_{23}		r_{2n_p}
	r_{31}	r_{32}	r_{33}	...	r_{3n_p}
	r_{ij}	...
User n_u	$r_{n_u 1}$	$r_{n_u 2}$	$r_{n_u 3}$		$r_{n_u n_p}$



Each column represents the preference for a single product across all users

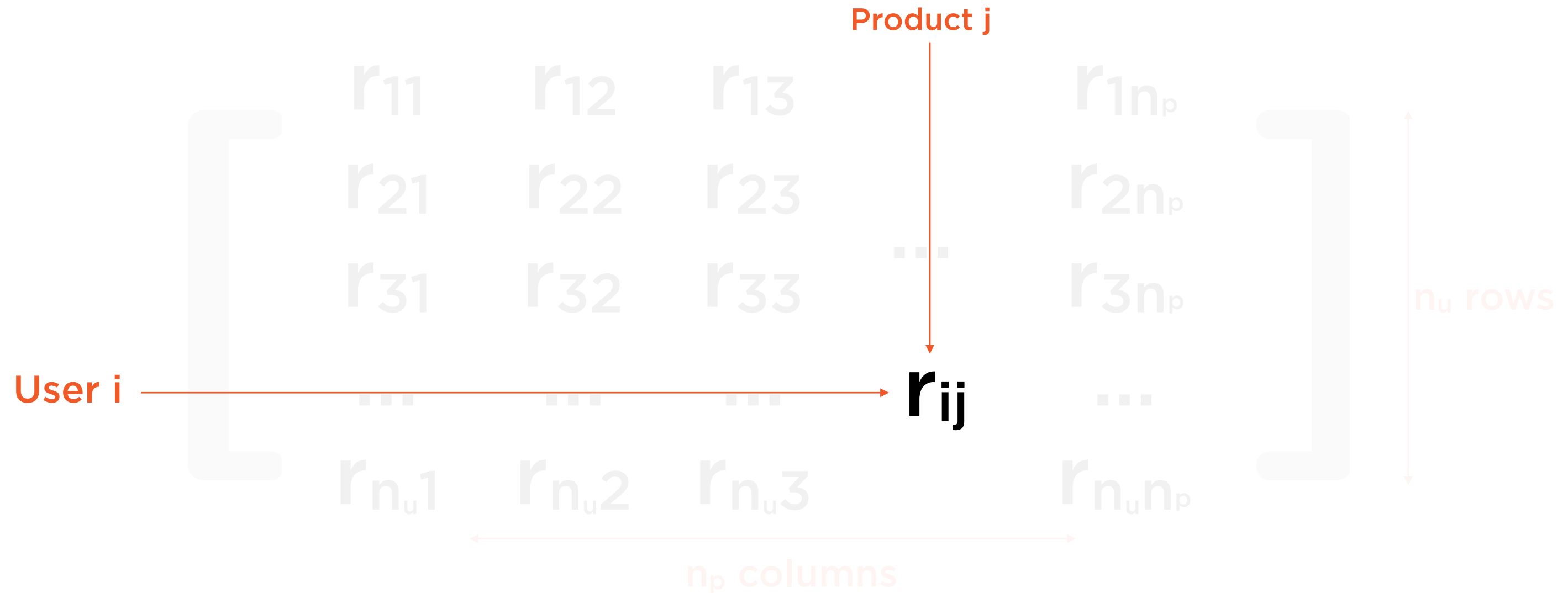
Ratings Matrix

					Product n_p
User 1	r_{11}	r_{12}	r_{13}		r_{1n_p}
User 2	r_{21}	r_{22}	r_{23}		r_{2n_p}
	r_{31}	r_{32}	r_{33}	...	r_{3n_p}
	r_{ij}	...
User n_u	r_{n_u1}	r_{n_u2}	r_{n_u3}		$r_{n_un_p}$



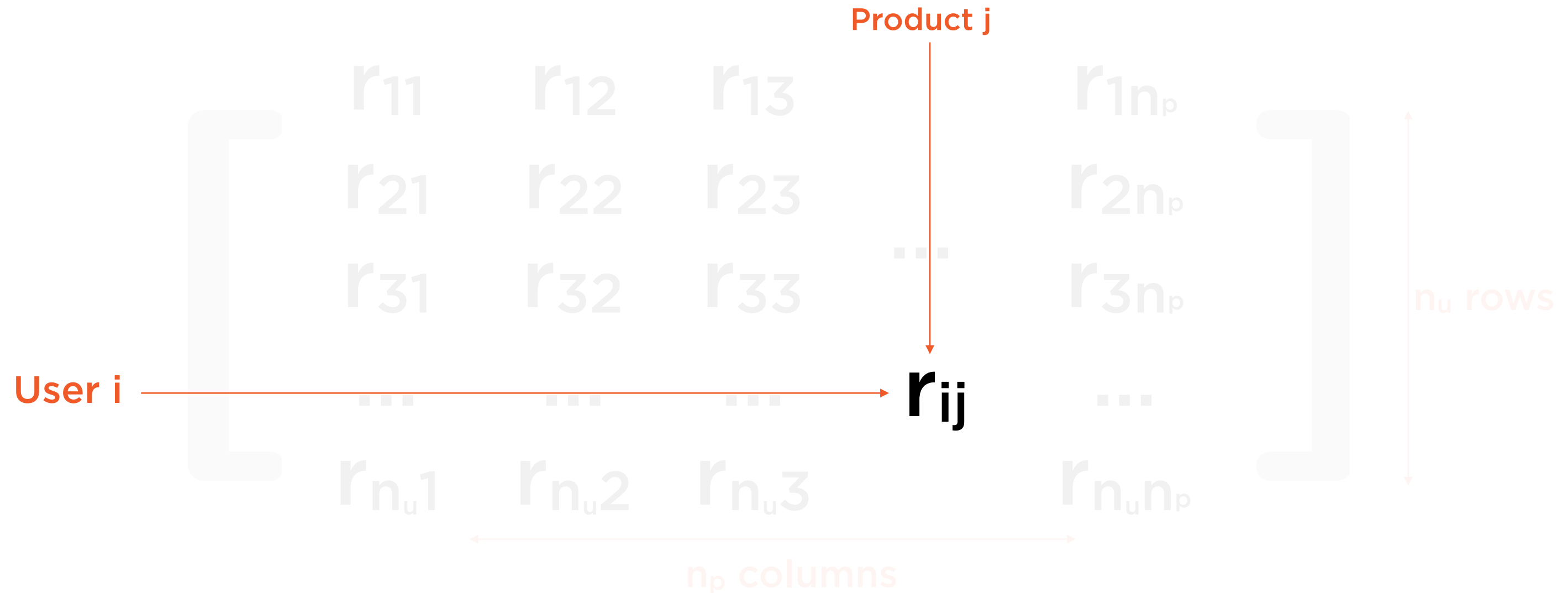
Each column represents the preference for a single product across all users

Ratings Matrix



Consider the rating of user i for product j

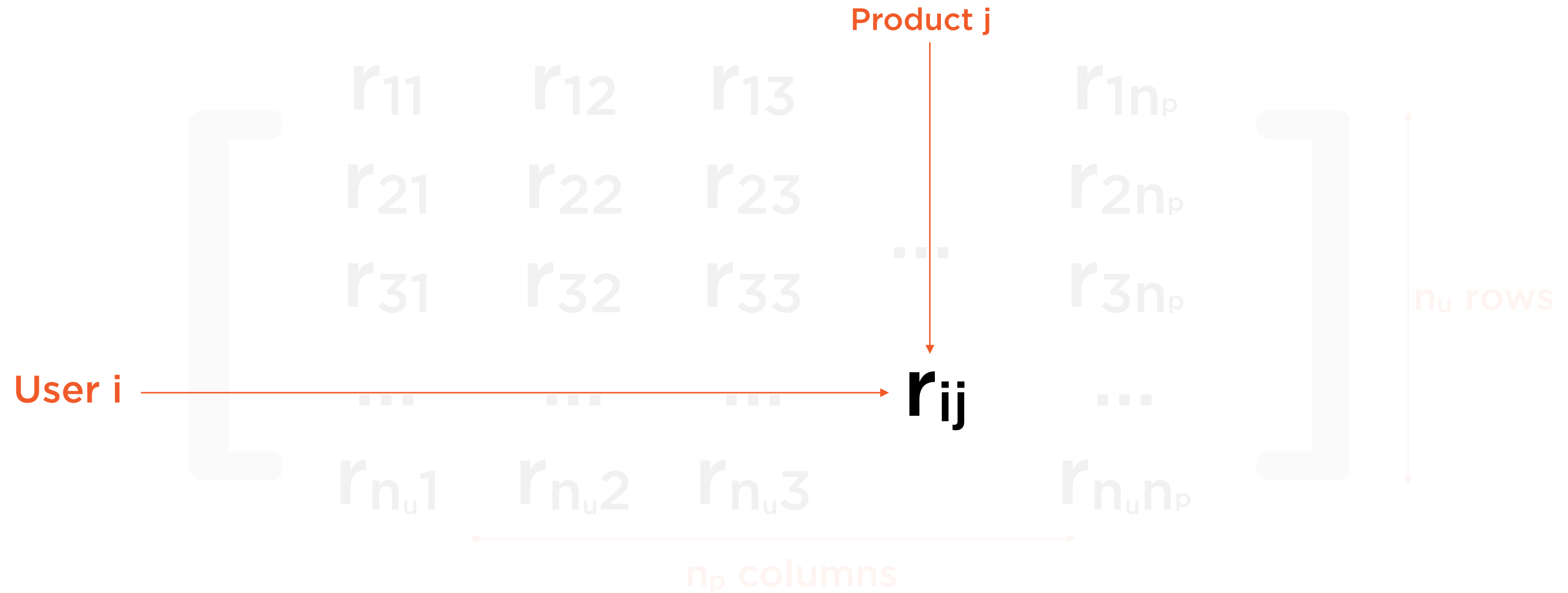
Ratings Matrix



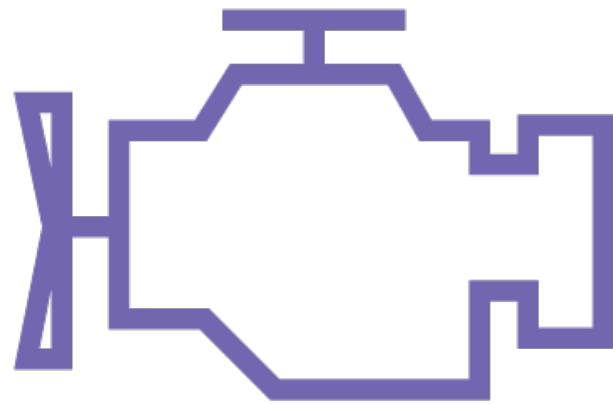
Very rarely, this user might actually have rated this product (e.g. by adding a rating + review)



Ratings Matrix



But usually, this value is initially missing and must be
estimated



Estimating Ratings Matrix

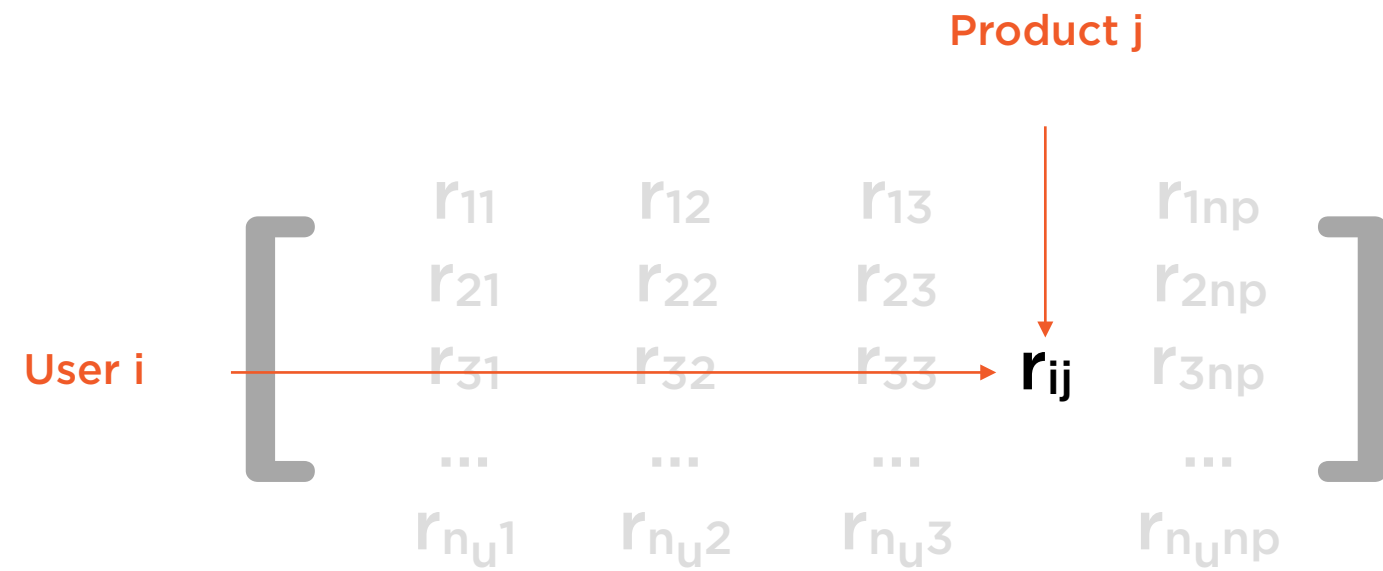
What if we could identify hidden factors that define this value?

This is a common technique called **latent factor analysis**

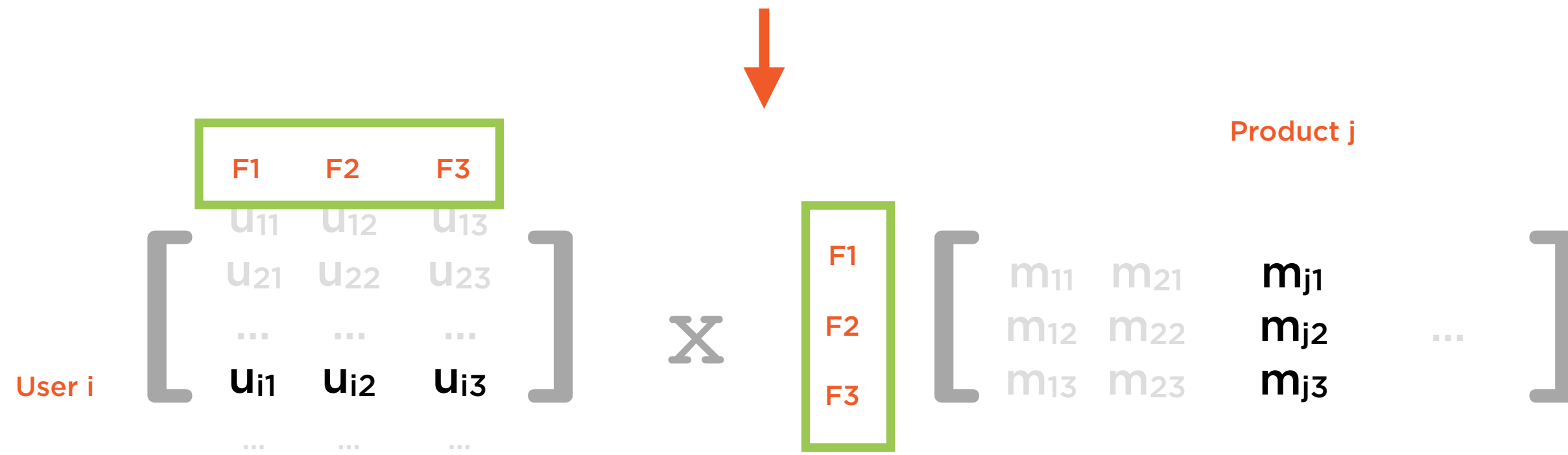
Pick a number of latent factors, say 3

$$n_f = 3$$

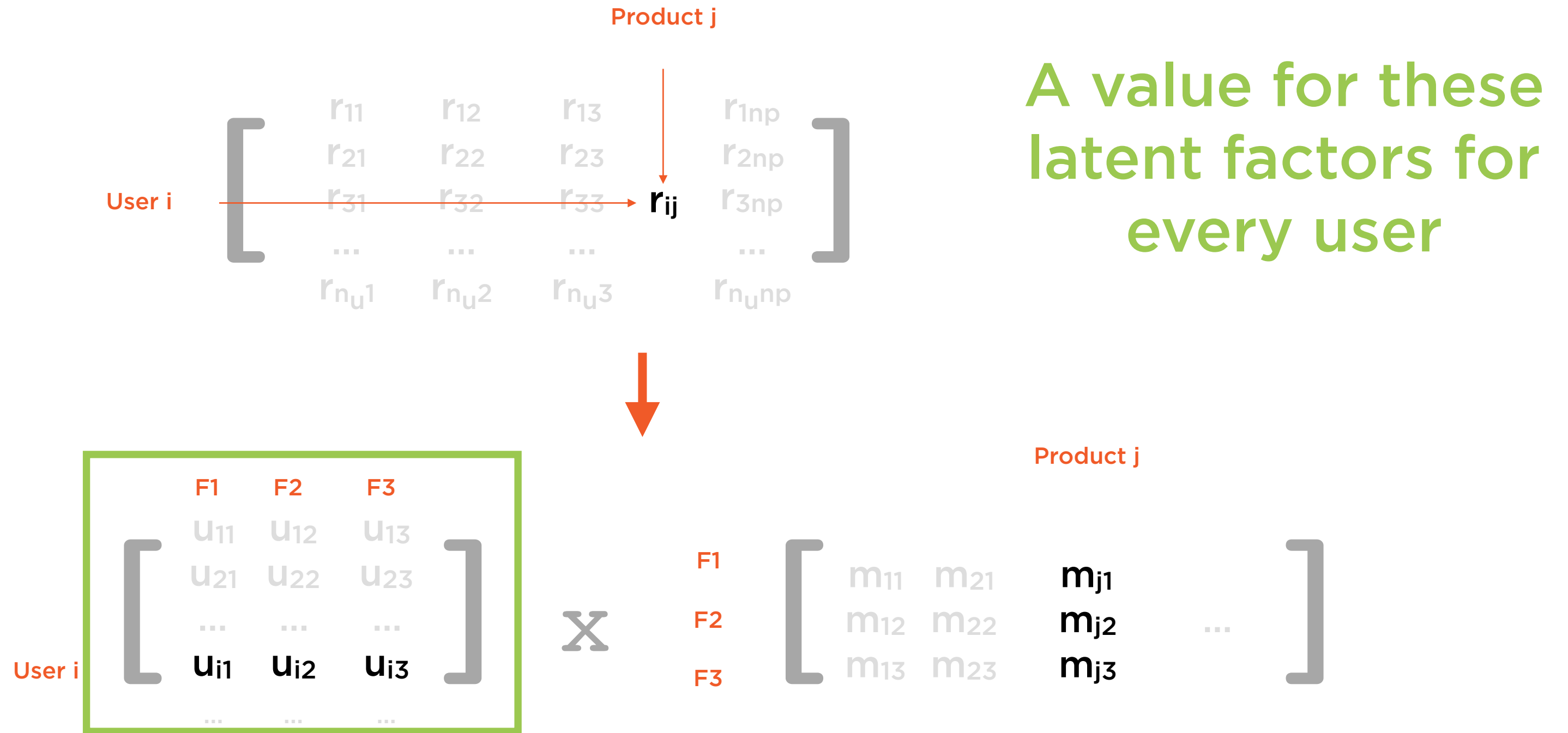
Ratings Matrix



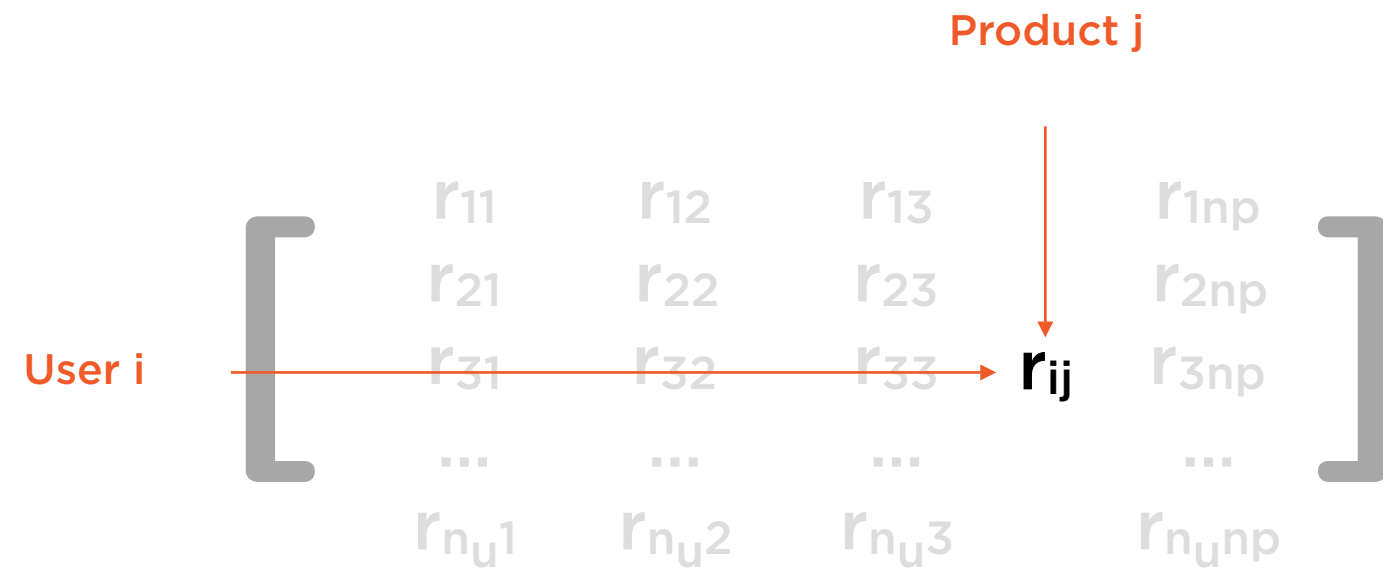
The 3 latent factors



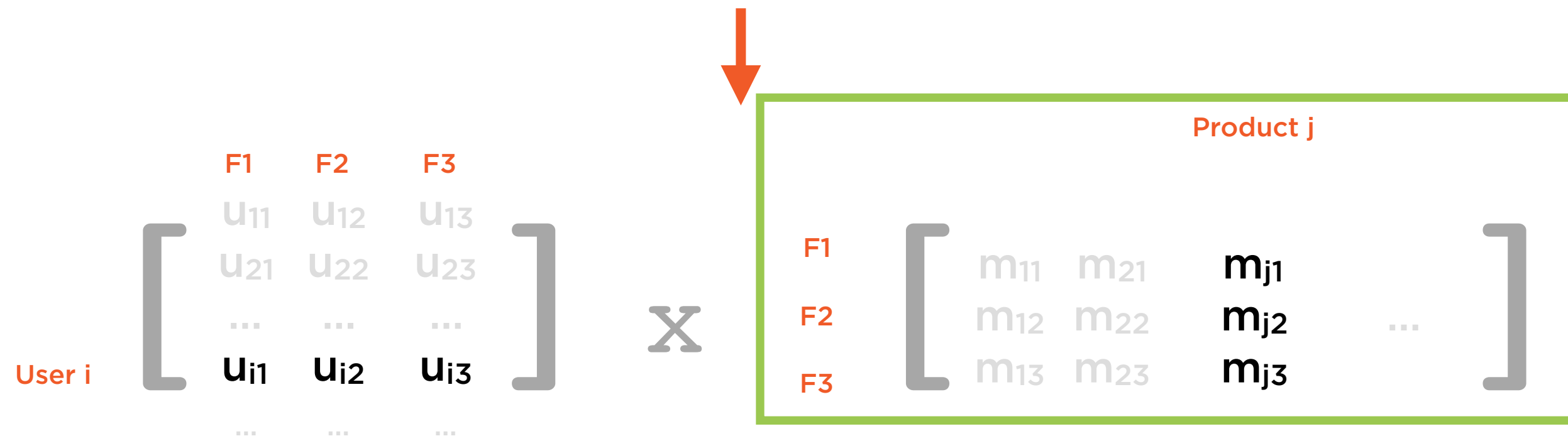
Ratings Matrix



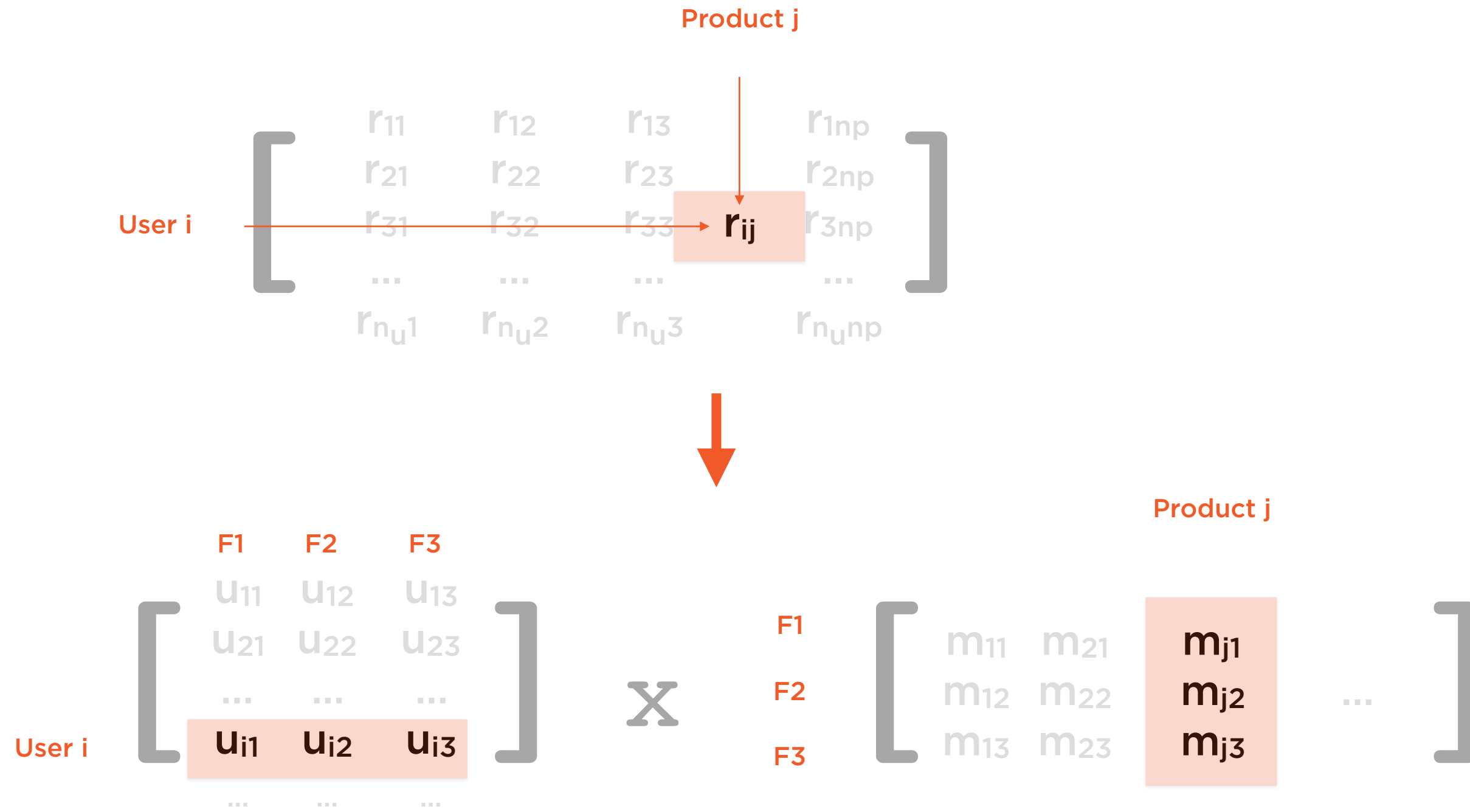
Ratings Matrix



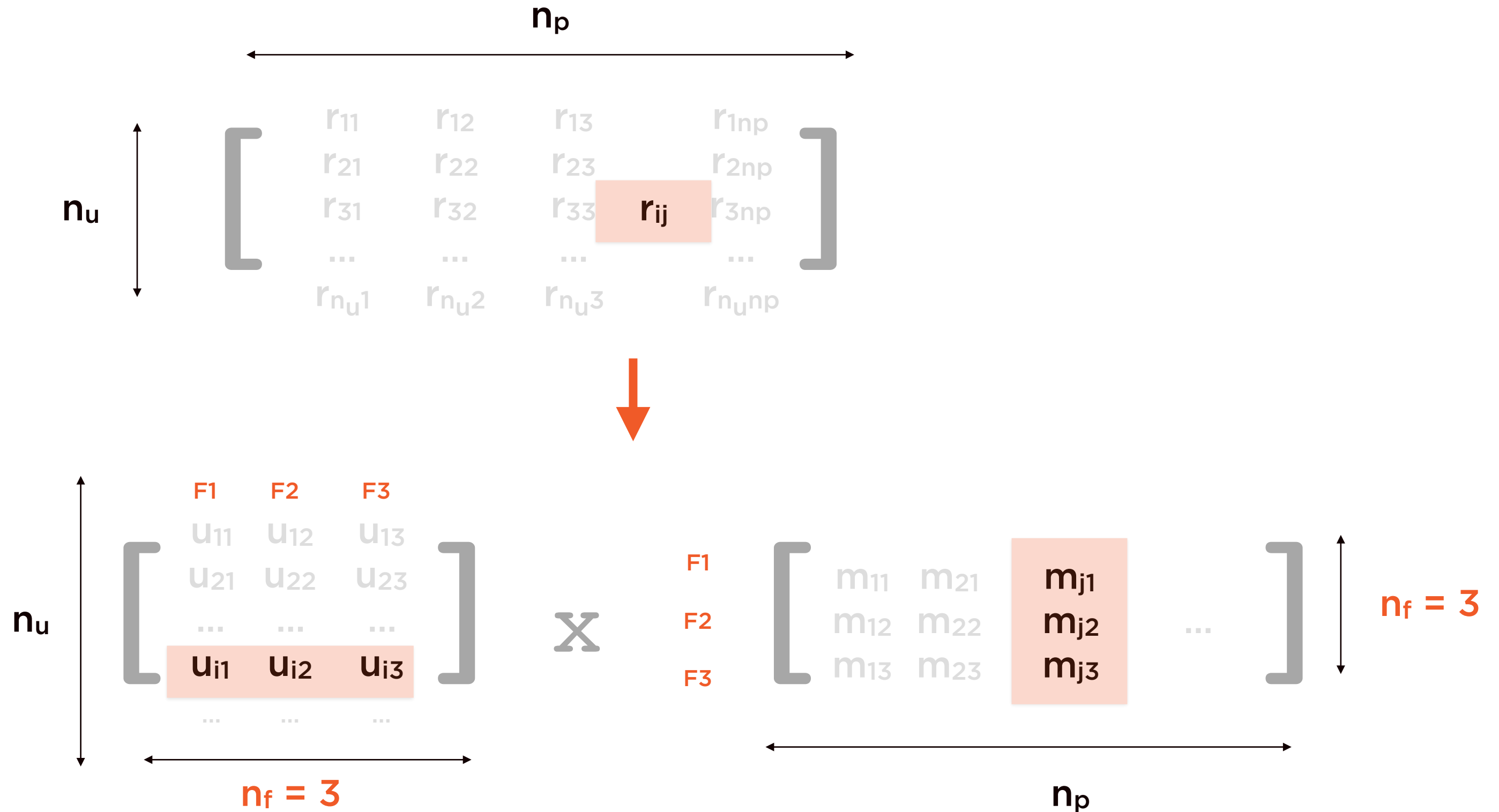
A value for these latent factors for every product



Ratings Matrix



Ratings Matrix



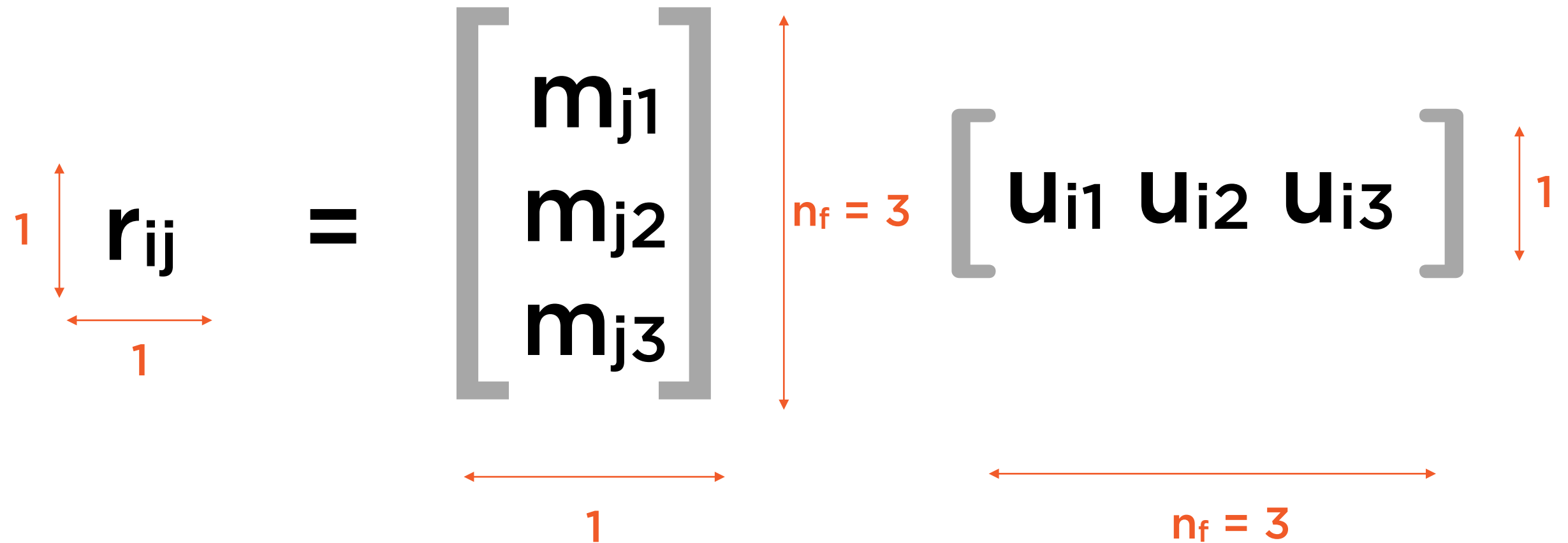
Matrix Factorization

The diagram illustrates the matrix factorization equation $r_{ij} = \begin{bmatrix} m_{j1} \\ m_{j2} \\ m_{j3} \end{bmatrix} \cdot \begin{bmatrix} u_{i1} & u_{i2} & u_{i3} \end{bmatrix}$. Dimensions are indicated by red arrows: r_{ij} is 1x1, the movie matrix is 3x1, the user matrix is 1x3, and the product is 1x1. The number of features $n_f = 3$ is shown for both the movie and user matrices.

$$\begin{matrix} \begin{matrix} \updownarrow 1 \\ r_{ij} \\ \leftarrow 1 \end{matrix} & = & \begin{matrix} \begin{bmatrix} m_{j1} \\ m_{j2} \\ m_{j3} \end{bmatrix} \\ \updownarrow n_f = 3 \\ \leftarrow 1 \end{matrix} & \cdot & \begin{matrix} \begin{bmatrix} u_{i1} & u_{i2} & u_{i3} \end{bmatrix} \\ \updownarrow 1 \\ \leftarrow n_f = 3 \end{matrix} \end{matrix}$$

Each entry in the user-rating matrix can be expressed as a matrix product

Matrix Factorization



The diagram illustrates the matrix factorization equation $r_{ij} = \begin{bmatrix} m_{j1} \\ m_{j2} \\ m_{j3} \end{bmatrix} \begin{bmatrix} u_{i1} & u_{i2} & u_{i3} \end{bmatrix}$. Dimensions are indicated by red arrows and labels: r_{ij} is 1x1; the first matrix is 3x1 with $n_f = 3$ rows; the second matrix is 1x3 with $n_f = 3$ columns. The value 1 is also shown as a dimension for the second matrix's columns.

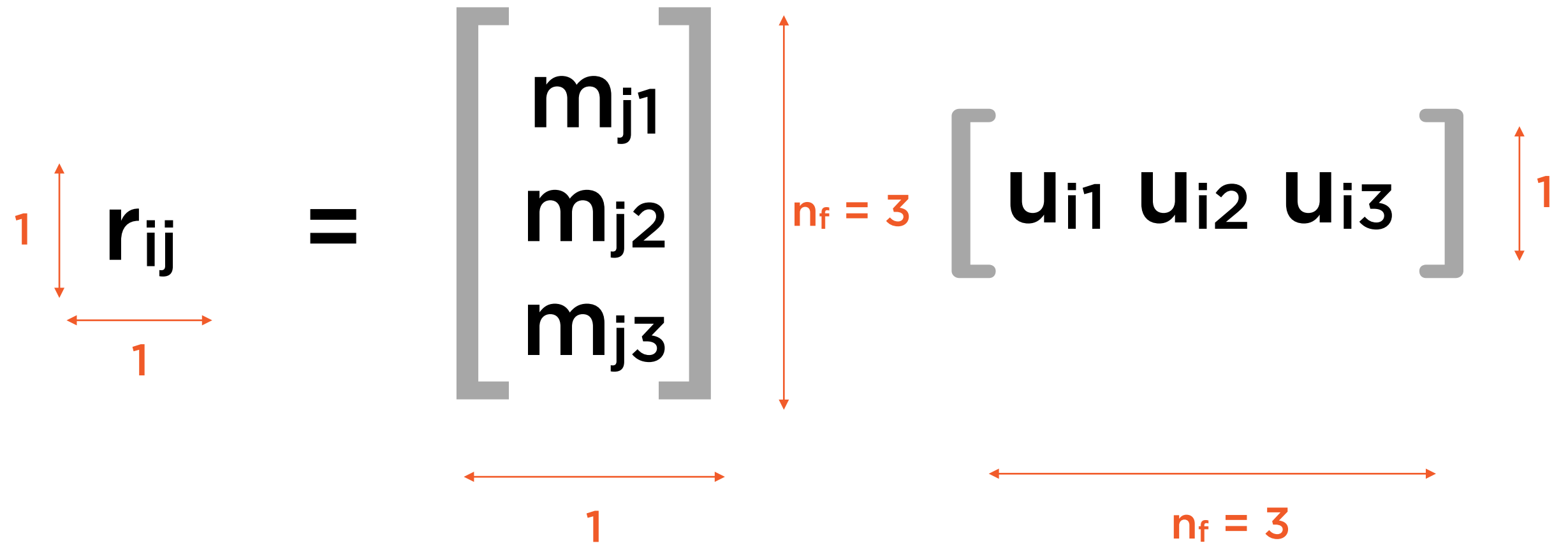
$$\begin{matrix} \text{1} \updownarrow \\ r_{ij} \\ \text{1} \leftarrow \end{matrix} = \begin{bmatrix} m_{j1} \\ m_{j2} \\ m_{j3} \end{bmatrix} \begin{matrix} \text{1} \updownarrow \\ \begin{bmatrix} u_{i1} & u_{i2} & u_{i3} \end{bmatrix} \\ \text{1} \leftarrow \end{matrix}$$

$n_f = 3$

$n_f = 3$

If we generalize this we get a system of linear equations to be solved

Matrix Factorization



The diagram illustrates the matrix factorization equation $r_{ij} = \begin{bmatrix} m_{j1} \\ m_{j2} \\ m_{j3} \end{bmatrix} \begin{bmatrix} u_{i1} & u_{i2} & u_{i3} \end{bmatrix}$. Dimensions are indicated by red arrows and labels: r_{ij} is 1x1; the first matrix is 3x1 with $n_f = 3$ rows; the second matrix is 1x3 with $n_f = 3$ columns. The variable n_f represents the number of features.

$$\begin{matrix} \text{1} \updownarrow \\ r_{ij} \\ \text{1} \leftarrow \end{matrix} = \begin{matrix} \begin{bmatrix} m_{j1} \\ m_{j2} \\ m_{j3} \end{bmatrix} \\ \text{1} \leftarrow \end{matrix} \begin{matrix} \begin{bmatrix} u_{i1} & u_{i2} & u_{i3} \end{bmatrix} \\ \text{1} \updownarrow \\ \text{1} \leftarrow \end{matrix}$$

$n_f = 3$

$n_f = 3$

Solving all of them simultaneously would allow us to estimate the entire matrix R

Matrix Factorization

The diagram illustrates the structure of a matrix R used in matrix factorization. On the left, the matrix R is shown with dimensions n_u (rows) and n_p (columns) indicated by red arrows. This is followed by an equals sign and a large gray bracketed matrix. The elements of this matrix are arranged in rows and columns, with the first row containing $r_{11}, r_{12}, r_{13}, \dots, r_{1n_p}$ and the first column containing $r_{21}, r_{31}, \dots, r_{n_u1}$. A general element r_{ij} is shown in the center. The bottom row contains $r_{n_u1}, r_{n_u2}, r_{n_u3}, \dots, r_{n_un_p}$. Red arrows indicate the dimensions: n_u rows (vertical arrow on the right) and n_p columns (horizontal arrow at the bottom).

$$R = \begin{bmatrix} r_{11} & r_{12} & r_{13} & \dots & r_{1n_p} \\ r_{21} & r_{22} & r_{23} & & r_{2n_p} \\ r_{31} & r_{32} & r_{33} & \dots & r_{3n_p} \\ \dots & \dots & \dots & r_{ij} & \dots \\ r_{n_u1} & r_{n_u2} & r_{n_u3} & & r_{n_un_p} \end{bmatrix}$$

n_u rows

n_p columns

Express this matrix as the product of two matrices, U and M

$$\begin{matrix} R & = & U & \times & M \\ \begin{matrix} n_u \text{ rows,} \\ n_p \text{ columns} \end{matrix} & & \begin{matrix} n_u \text{ rows,} \\ n_f \text{ columns} \end{matrix} & & \begin{matrix} n_f \text{ rows,} \\ n_p \text{ columns} \end{matrix} \end{matrix}$$

n_f is a hyperparameter

Estimating Rating Matrix

$$R = U \times M$$

n_u rows,
 n_p columns n_u rows,
 n_f columns n_f rows,
 n_p columns

n_f is a hyperparameter

“rank”

“Number of latent factors”

“Dimensionality of feature space”

Estimating Rating Matrix

$$R = U \times M$$

n_u rows,
 n_p columns n_u rows,
 n_f columns n_f rows,
 n_p columns

If R were available...

...Many matrix techniques to find U,M

e.g. Singular Value Decomposition

(Used in PCA)

$$\underset{\substack{n_u \text{ rows,} \\ n_p \text{ columns}}}{R} = \underset{\substack{n_u \text{ rows,} \\ n_f \text{ columns}}}{U} \times \underset{\substack{n_f \text{ rows,} \\ n_p \text{ columns}}}{M}$$

Estimating Rating Matrix

But R is not available and needs to be estimated

Use **Alternating-Least-Squares (ALS)**

Standard numerical algorithm

Alternating Least Squares (ALS)

Minimize

$$\sum_{i,j} (r_{ij} - u_i m_j)^2$$

To find

U, M

The value of U and M define the “best” rating matrix

$$\mathbf{R} = \mathbf{U} \times \mathbf{M}$$

Step 1: Initialize M

- ◀ Assign average rating for that product as first row
 - ◀ Small random numbers for other rows

Step 2:
Fix M , solve to find U

- ◀ Solve to minimize squared errors

Step 3:
Fix U , solve to find M

- ◀ Solve to minimize squared errors

Step 4:
If stopping criterion not met
Repeat Steps 2 and 3

- ◀ Stop if RMSE on training data lower than some threshold

$$R = U \times M$$

n_u rows,
 n_p columns n_u rows,
 n_f columns n_f rows,
 n_p columns

Estimating Rating Matrix

Each element of U, M is a free parameter

The number of free parameters is very large

Likely to lead to overfitting

Add regularization to penalize large parameters

ALS-WR

$$R = U \times M$$

n_u rows,
 n_p columns n_u rows,
 n_f columns n_f rows,
 n_p columns

Alternating-Least-Squares (ALS)

Weighted Regularization (WR)

ALS-WR

Minimize

$$\sum_{i,j} (r_{ij} - u_i m_j)^2 + \lambda \left(\sum_i n_u^i u_i^2 + \sum_j n_m^j m_j^2 \right)$$

To find

U, M

λ is a hyperparameter that penalizes complex models

ALS-WR

Minimize

$$\sum_{i,j} (r_{ij} - u_i m_j)^2 + \lambda \left(\sum_i n_u^i u_i^2 + \sum_j n_m^j m_j^2 \right)$$

To find

U, M

λ is a hyperparameter that penalizes complex models

Difficulties in Estimation

Sparsity

Most initial entries will be missing and need estimation

Cold Start

New users or products with no history

Computational Intensity

Millions of users, millions of products

Difficulties in Estimation

Sparsity

Most initial entries will be missing and need estimation

Cold Start

New users or products with no history

Computational Intensity

Millions of users, millions of products

Sparsity

Most initial entries will be missing and need estimation

Choice of Rating Measure

Smart choice of rating measure can help mitigate

Two types

- explicit
- implicit

Spark supports both

Choice of Rating Measure



Implicit

Number of times viewed, length
of time listened...



Explicit

Reviews or ratings manually
entered, explicit actions
performed...

Explicit and Implicit Feedback

Explicit Feedback

Star ratings, thumbs-up/down buttons
clicked

Far less available

Easy to use directly in recommendation
algorithms

Explicit ratings express preference
between items

Can be used to order preferences

Implicit Feedback

Browsing history, mouse movements,
time spent listening/viewing

Easily available and plentiful

Hard to use directly in recommendation
algorithms for several reasons

Implicit ratings do not express
preferences between items

Can only use to model confidence we
have in a particular observation



Challenges of Implicit Feedback

No negative feedback

Inherently noisy (viewing to buy a gift?)

Unlike explicit feedback - no preference or order

However implicit feedback does give us confidence in data items



Spark's Approach to Implicit Feedback

Treat elements of R as explicit feedback

But do so in a smart, efficient way

Allows use of implicit feedback without compromising correctness

Basically - model R as strength of observations



Spark's Approach to Implicit Feedback

**“Collaborative Filtering for Implicit Feedback
Datasets”**

- Yifan Hu, Yehuda Koren, Chris Volinsky (2008)



Spark's Approach to Implicit Feedback

Optimization: ALS-WR can be used

- RMSE fine for use here

Evaluation: RMSE will not suffice

- RMSE won't work (no negative feedback)
- Use rank measure instead
- Math is complex and out-of-scope

Difficulties in Estimation

Sparsity

Most initial entries will be missing and need estimation

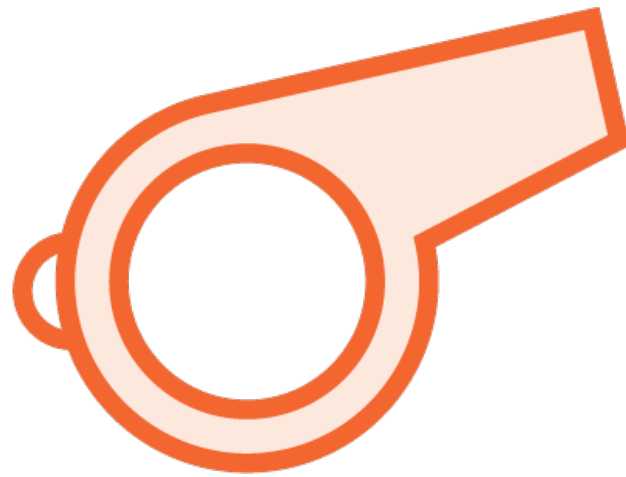
Cold Start

New users or products with no history

Computational Intensity

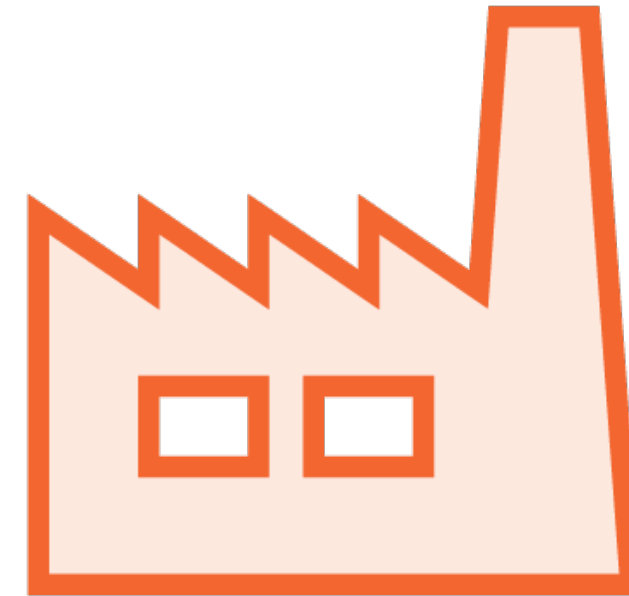
Millions of users, millions of products

Cold Start



During Validation

Encounter user or product not in the training dataset



In Production

Entirely new products or users enter the system

Cold Start

New users or products
with no history

Cold Start Strategies

Different strategies needed for these two cases

- **In production: Spark assigns NaN (default)**
- **In cross-validation: Spark will “drop” row**

“nan” and “drop” are two permissible strategies

Difficulties in Estimation

Sparsity

Most initial entries will be missing and need estimation

Cold Start

New users or products with no history

Computational Intensity

Millions of users, millions of products



Spark's Approach to Scalability

“Collaborative Filtering for Implicit Feedback Datasets”

- Yifan Hu, Yehuda Koren, Chris Volinsky (2008)

This algorithm not only solves for implicit feedback...

...it also scales linearly with dataset size!

Demo

Recommendation systems using explicit ratings in spark.ml

Demo

Recommendation systems using implicit ratings in spark.ml

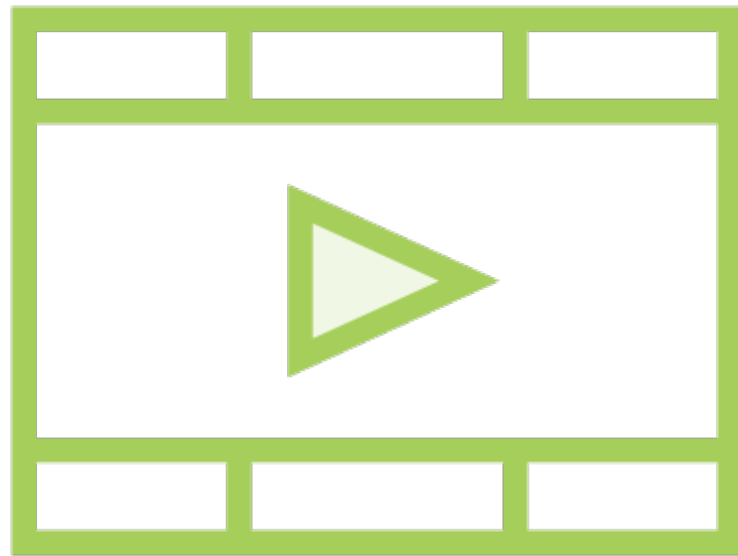
Overview

Collaborative filtering algorithms are used to make product recommendations to users

Spark offers estimators which use the ALS method to implement collaborative filtering

Recommendations can be based on explicit and implicit ratings

Related Courses - Spark



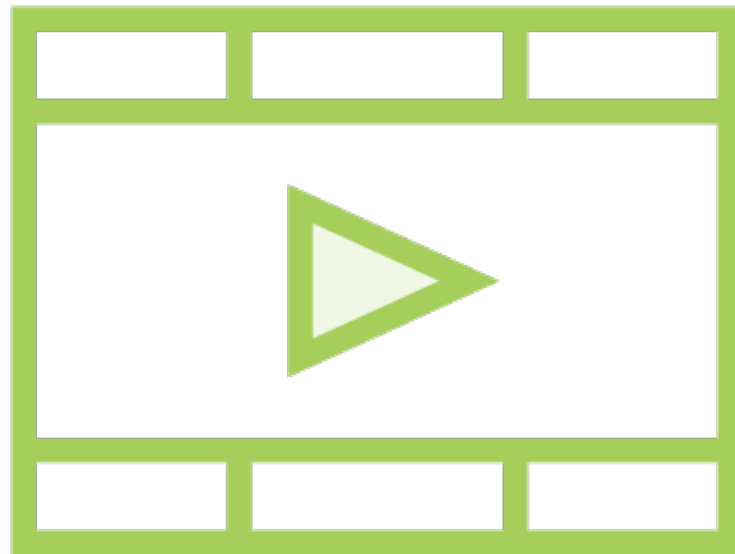
Handling Fast Data with Apache Spark SQL and Streaming

- Programming in Spark 2 using Scala

Getting Started with Stream Processing with Spark Streaming

- Stream processing with Spark 1.x in Python

Related Courses - ML



Understanding the Foundations of TensorFlow

- Introduction to TensorFlow to build NNs

TensorFlow: Getting Started