# Dealing with sparsity

## BUILDING RECOMMENDATION ENGINES IN PYTHON
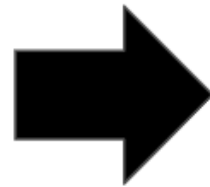
**Rob O'Callaghan**
Director of Data

# Sparse matrices

# Sparse matrices

# Sparse matrices



Sparsity = $\dfrac{\text{Empty Values}}{\text{Total Cells}}$

# Measuring sparsity

```python
print(book_rating_df)
```

| title | The Great Gatsby | The Catcher in the Rye | Fifty Shades of Grey |
|---|---|---|---|
| User | | | |
| User_233 | 3.0 | NaN | NaN |
| User_651 | NaN | 5.0 | 4.0 |
| User_965 | 4.0 | 3.0 | NaN |
| ... | ... | ... | ... |

# Measuring sparsity

```python
number_of_empty = book_ratings_df.isnull().values.sum()
total_number = user_ratings_df.size
sparsity = number_of_empty/total_number
print(sparsity)
```

```
0.0114
```

# Why sparsity matters

# Why sparsity matters

|  | Item 1 | Item 2 | Item 3 | Item 4 | Item 5 | Item 6 |
|---|---|---|---|---|---|---|
| User 1 |  | 1 |  |  | ? |  |
| User 2 |  |  |  |  | 5 |  |
| User 3 |  |  |  | 3 |  |  |
| User 4 |  |  |  |  |  | 4 |
| User 5 | 2 |  |  |  |  |  |
| User 6 |  |  | 5 |  | 1 |  |

# Why sparsity matters

# Why sparsity matters

# Measuring sparsity per column

```python
user_ratings_df.notnull().sum()
```

```
The Pelican Brief             1
Snow Crash                    1
The Great Gatsby             12
Fifty Shades of Grey          9
Leviathan                     1
                             ..
```

# Matrix factorization

| | Item 1 | Item 2 | Item 3 | Item 4 | Item 5 | Item 6 |
|--------|--------|--------|--------|--------|--------|--------|
| User 1 | ? | 1 | ? | ? | ? | ? |
| User 2 | ? | ? | ? | ? | 5 | ? |
| User 3 | ? | ? | ? | 3 | ? | ? |
| User 4 | ? | ? | ? | ? | ? | 4 |
| User 5 | 2 | ? | ? | ? | ? | ? |
| User 6 | ? | ? | 5 | ? | 1 | ? |

**Original DataFrame**

# Matrix factorization



Original DataFrame

DataFrame Factors

# Matrix factorization



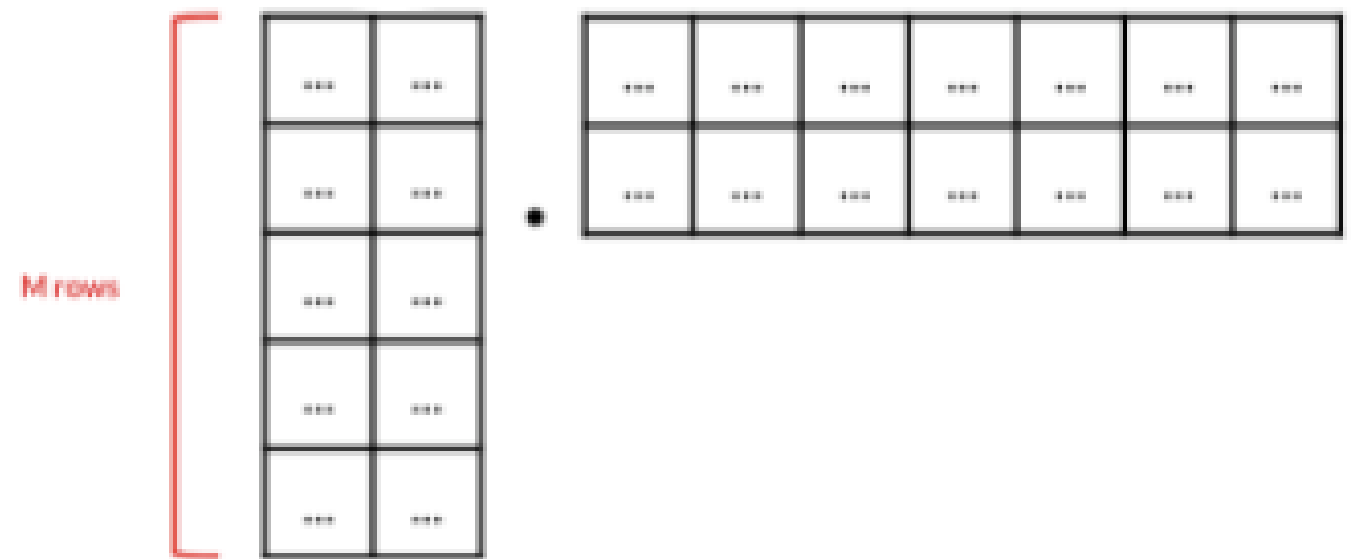Original DataFrame      DataFrame Factors      Filled DataFrame
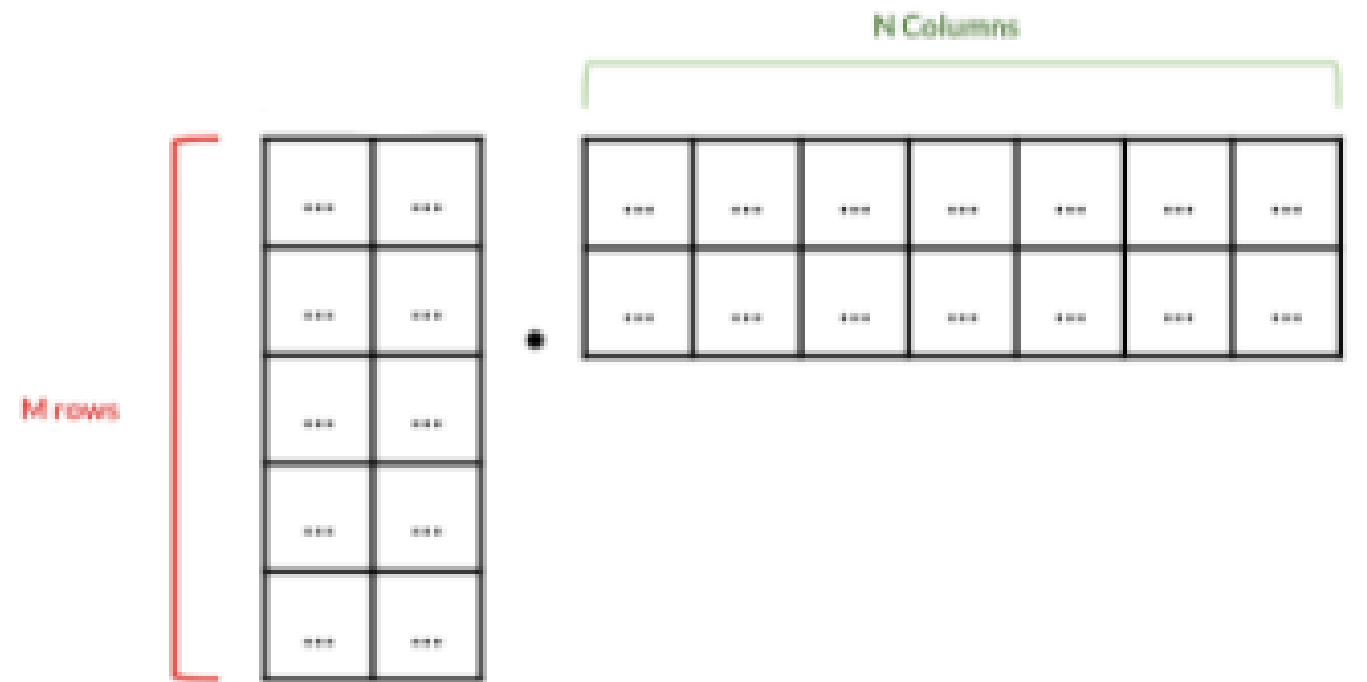
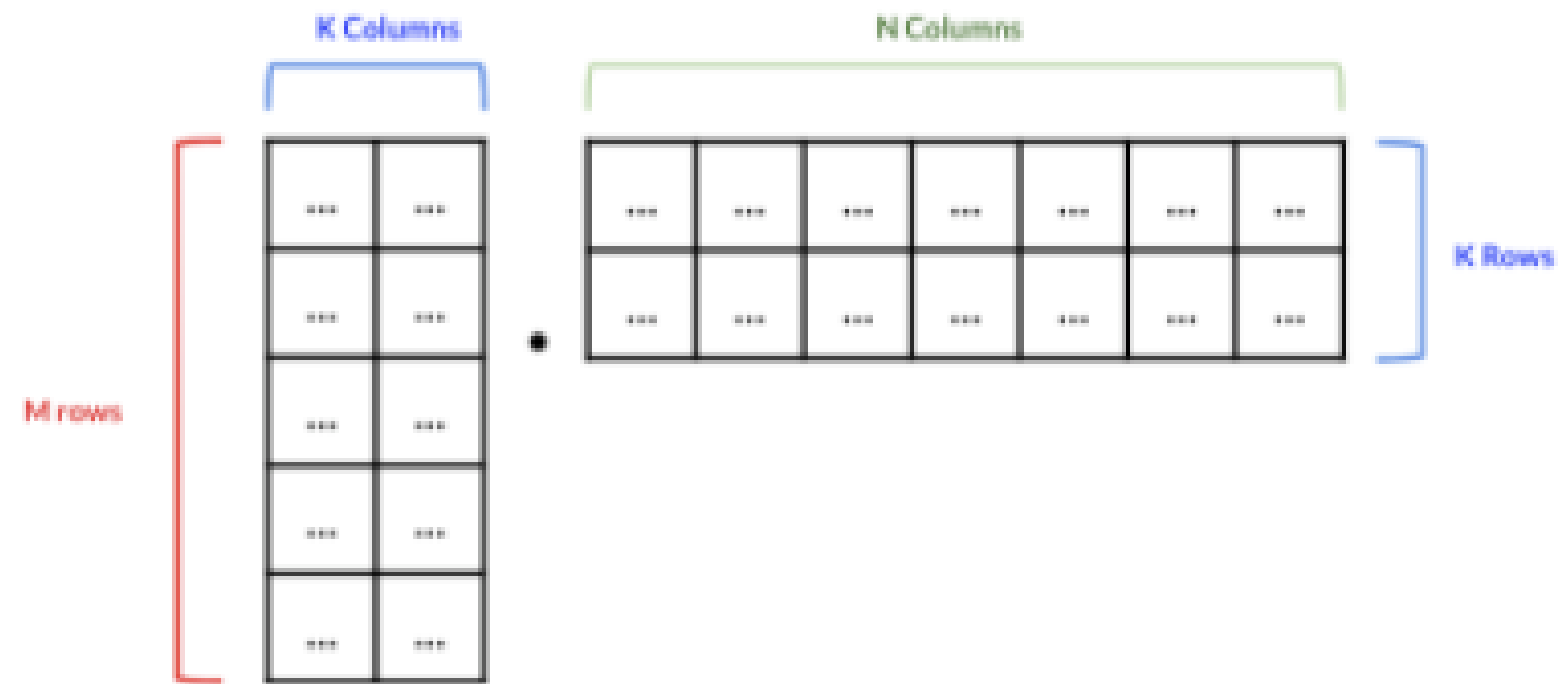# Matrix multiplication

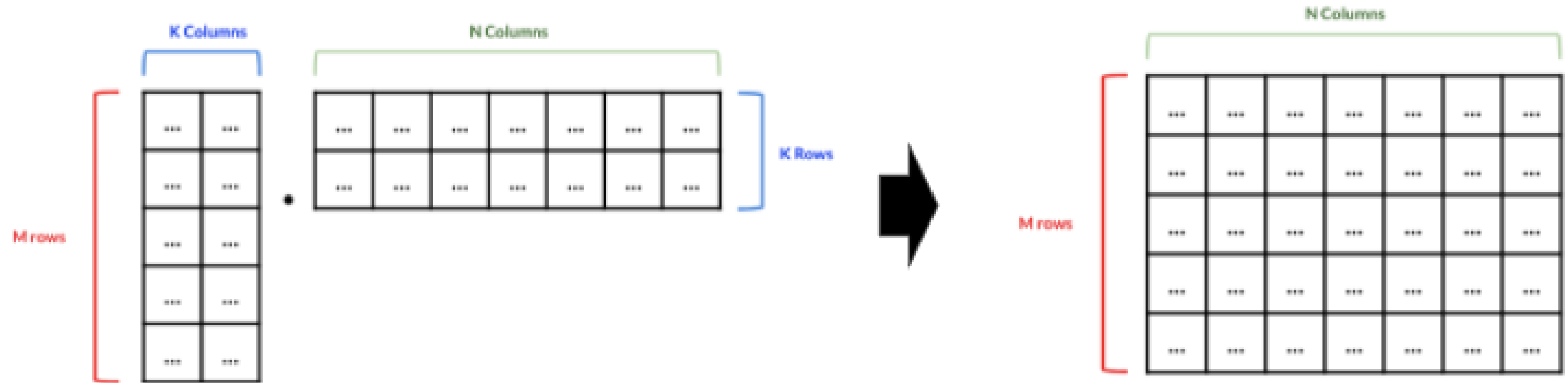# Matrix multiplication

# Matrix multiplication

# Matrix multiplication

# Matrix multiplication

# Matrix multiplication

```python
print(matrix_x)
```

```
[[4, 1],
 [2, 2],
 [3, 3]]
```

```python
print(matrix_b)
```

```
[[1, 0, 4],
 [0, 1, 6]]
```

# Matrix multiplication

```python
import numpy as np

dot_product = np.dot(matrix_x, matrix_b)
print(dot_product)
```

```
[[ 4  1 22]
 [ 2  2 20]
 [ 3  3 30]]
```
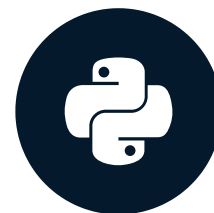
# Let's practice!

datacamp

# Why this helps with sparse matrices



Original DataFrame

# Why this helps with sparse matrices



Original DataFrame

DataFrame Factors

# Why this helps with sparse matrices



Original DataFrame

DataFrame Factors

Filled DataFrame

# What matrix factorization looks like

# What matrix factorization looks like

# What matrix factorization looks like

# What matrix factorization looks like

**BUILDING RECOMMENDATION ENGINES IN PYTHON**

# Latent features

# Latent features

# Latent features

# Information loss



Original Matrix

# Information loss



Original Matrix          Matrix Factors

# Information loss

Original Matrix

Matrix Factors

Output Matrix

# Information loss

**Original Matrix**

**Matrix Factors**

**Output Matrix**

Similar but not
identical

# Let's practice!

BUILDING RECOMMENDATION ENGINES IN PYTHON

# Singular value decomposition (SVD)

## BUILDING RECOMMENDATION ENGINES IN PYTHON

**Rob O'Callaghan**
Director of Data

# What SVD does



Original Data

| | 5 | |
|---|---|---|
| | | 3 |
| 4 | | |
| | | 2 |

$U$

| 0.8 | 0.2 |
|---|---|
| 0.5 | 0.4 |
| 0.7 | 0.1 |
| 0.1 | 0.9 |

$\Sigma$

| 11 | 0 |
|---|---|
| 0 | 2.5 |

$V^t$

| ... | ... | ... |
|---|---|---|
| ... | ... | ... |

# What SVD does

**Original Data**

| | | |
|---|---|---|
| | 5 | |
| | | 3 |
| 4 | | |
| | | 2 |

➡️

**U**

| | |
|---|---|
| 0.8 | 0.2 |
| 0.5 | 0.4 |
| 0.7 | 0.1 |
| 0.1 | 0.9 |

**Σ**

| | |
|---|---|
| 11 | 0 |
| 0 | 2.5 |

**V<sup>t</sup>**

| | | |
|---|---|---|
| ... | ... | ... |
| ... | ... | ... |

# What SVD does

# What SVD does



Original Data

| | 5 | |
|---|---|---|
| | | 3 |
| 4 | | |
| | | 2 |

$U$

| 0.8 | 0.2 |
|---|---|
| 0.5 | 0.4 |
| 0.7 | 0.1 |
| 0.1 | 0.9 |

$\Sigma$

| 11 | 0 |
|---|---|
| 0 | 2.5 |

$V^t$

| ... | ... | ... |
|---|---|---|
| ... | ... | ... |

# Prepping our data

```python
print(book_ratings_df.shape)
```

```
(220, 500)
```

```python
avg_ratings = book_ratings_df.mean(axis=1)
print(avg_ratings)
```

```
array([[4.5 ],
       [3.5],
       [2.5],
       [3.5],
        ...
       [2.2]])
```

# Prepping our data

```python
user_ratings_pivot_centered = user_ratings_df.sub(avg_ratings, axis=0)
user_ratings_df.fillna(0, inplace=True)
print(user_ratings_df)
```

```
             The Great Gatsby    The Catcher in the Rye    Fifty Shades of Grey
User_233                 0.0                       0.0                     0.0
User_651                 0.0                       0.5                    -0.5
User_965                 0.5                      -0.5                     0.0
    ...                  ...                       ...                     ...
```

# Applying SVD

```python
from scipy.sparse.linalg import svds
U, sigma, Vt = svds(user_ratings_pivot_centered)
```

```python
print(U.shape)
```

```
(610, 6)
```

```python
print(Vt.shape)
```

```
(6, 1000)
```

# Applying SVD

```python
print(sigma)
```

```
[3.0, 4.8, -12.6, -3.8, 8.2, 7.3]
```

```python
sigma = np.diag(sigma)
print(sigma)
```

```
array([   3.0   ,   0.    ,   0.    ,   0.    ,   0.    ,   0.    ],
      [  0.    ,   4.8   ,   0.    ,   0.    ,   0.    ,   0.    ],
      [  0.    ,   0.    , -12.6   ,   0.    ,   0.    ,   0.    ],
      [  0.    ,   0.    ,   0.    ,  -3.8   ,   0.    ,   0.    ],
      [  0.    ,   0.    ,   0.    ,   0.    ,   8.2   ,   0.    ],
      [  0.    ,   0.    ,   0.    ,   0.    ,   0.    ,   7.3   ]),
```

# Getting the final matrix



$U$  $\Sigma$  $V^t$  $\Rightarrow$  $U\Sigma V^t$

# Getting the final matrix

# Getting the final matrix

# Getting the final matrix

# Calculating the product in Python

```
recalculated_ratings =          np.dot(U, sigma)
```

# Calculating the product in Python

```python
recalculated_ratings = np.dot(np.dot(U, sigma), Vt)
print(recalculated_ratings)
```

```
[[  0.1      -0.9      -3.6.     ...   ]
 [ -2.3       0.5      -0.5      ...   ]
 [  0.5      -0.5       2.0      ...   ]
 [ ...       ...       ...       ...   ]]
```

# Add averages back

```python
recalculated_ratings = recalculated_ratings + avg_ratings.values.reshape(-1, 1)
print(recalculated_ratings)
```

```
[[ 4.6      3.6      0.9      ...      ]
 [ 1.8      4.0      3.0      ...      ]
 [ 3.0      2.0      4.5      ...      ]
 [ ...      ...      ...      ...      ]]
```

```python
print(book_ratings_df)
```

```
[[ 5.0      4.0       NA      ...      ]
 [  NA      4.0      3.0      ...      ]
 [ 3.0      2.0       NA      ...      ]
 [ ...      ...      ...      ...      ]]
```

# Let's practice!

datacamp

# Validating your predictions

BUILDING RECOMMENDATION ENGINES IN PYTHON

**Rob O'Callaghan**
Director of Data

# Hold-out sets



Most Machine Learning
Models

Recommendation Engines

# Hold-out sets

|  | Item 1 | Item 2 | Item 3 | Item 4 | Item 5 | Item 6 |
|---|---|---|---|---|---|---|
| User 1 |  |  |  |  |  |  |
| User 2 |  |  |  |  |  |  |
| User 3 |  |  |  |  |  |  |
| User 4 |  |  |  |  |  |  |
| User 5 |  |  |  |  |  |  |
| User 6 |  |  |  |  |  |  |

**Most Machine Learning Models**

|  | Item 1 | Item 2 | Item 3 | Item 4 | Item 5 | Item 6 |
|---|---|---|---|---|---|---|
| User 1 |  |  |  |  |  |  |
| User 2 |  |  |  |  |  |  |
| User 3 |  |  |  |  |  |  |
| User 4 |  |  |  |  |  |  |
| User 5 |  |  |  |  |  |  |
| User 6 |  |  |  |  |  |  |

**Recommendation Engines**

# Hold-out sets



Most Machine Learning
Models

Recommendation Engines

# Hold-out sets



Most Machine Learning Models

Recommendation Engines

# Hold-out sets



Most Machine Learning Models

Recommendation Engines

# Hold-out sets



Most Machine Learning Models

Recommendation Engines

# Separating the hold-out set

```
actual_values = act_ratings_df.iloc[:20, :100].values
act_ratings_df.iloc[:20, :100] = np.nan
```

Generate predictions as before.

```
predicted_values = calc_pred_ratings_df.iloc[:20, :100].values
```

# Masking the hold-out set

```python
mask = ~np.isnan(actual_values)
```

```python
print(actual_values[mask])
```

```
[4.  4.  5.  3.  3.  ...]
```

```python
print(predicted_values[mask])
```

```
[3.76, 4.35,  4.95,  3.5869079 3.686337   ...]
```

# Introducing RMSE (root mean squared error)

| Predicted | Actual |
|:---------:|:------:|
| 4 | 5 |
| 3 | 3 |
| 2 | 4 |

# Introducing RMSE (root mean squared error)

| Predicted | Actual | Difference |
|:---:|:---:|:---:|
| 4 | 5 | 1 |
| 3 | 3 | 0 |
| 2 | 4 | 2 |

# Introducing RMSE (root mean squared error)

| Predicted | Actual | Difference | Difference$^2$ |
|:---------:|:------:|:----------:|:--------------:|
| 4 | 5 | 1 | 1 |
| 3 | 3 | 0 | 0 |
| 2 | 4 | 2 | 4 |

# Introducing RMSE (root mean squared error)

| Predicted | Actual | Difference | Difference$^2$ |
|:---------:|:------:|:----------:|:--------------:|
| 4 | 5 | 1 | 1 |
| 3 | 3 | 0 | 0 |
| 2 | 4 | 2 | 4 |

$$\sqrt{\frac{\text{Sum of Diff}^2}{\text{Count}}}$$

# Introducing RMSE (root mean squared error)

| Predicted | Actual | Difference | Difference$^2$ |
|:---:|:---:|:---:|:---:|
| 4 | 5 | 1 | 1 |
| 3 | 3 | 0 | 0 |
| 2 | 4 | 2 | 4 |

$$\sqrt{\dfrac{5}{3}}$$

# Introducing RMSE (root mean squared error)

| Predicted | Actual | Difference | Difference$^2$ |
|-----------|--------|------------|----------------|
| 4 | 5 | 1 | 1 |
| 3 | 3 | 0 | 0 |
| 2 | 4 | 2 | 4 |

$$\sqrt{\frac{5}{3}} \implies 1.29$$

# RMSE in Python

```python
from sklearn.metrics import mean_squared_error

print(mean_squared_error(actual_values[mask],
                         predicted_values[mask],
                         squared=False))
```
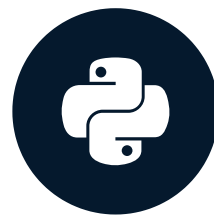
```
3.6223997
```

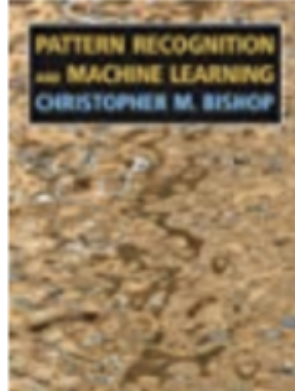# Let's practice!

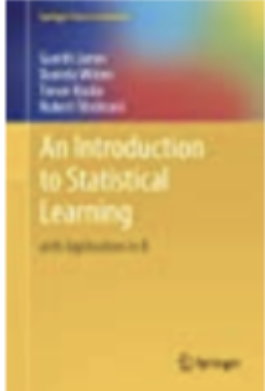BUILDING RECOMMENDATION ENGINES IN PYTHON
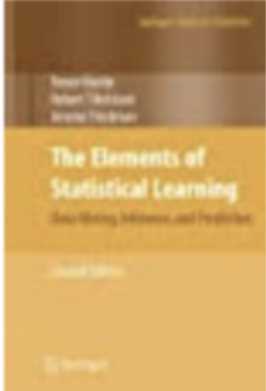
# Wrap up

## BUILDING RECOMMENDATION ENGINES IN PYTHON

**Rob O'Callaghan**
Director of Data

# Non-personalized models

## Frequently bought together
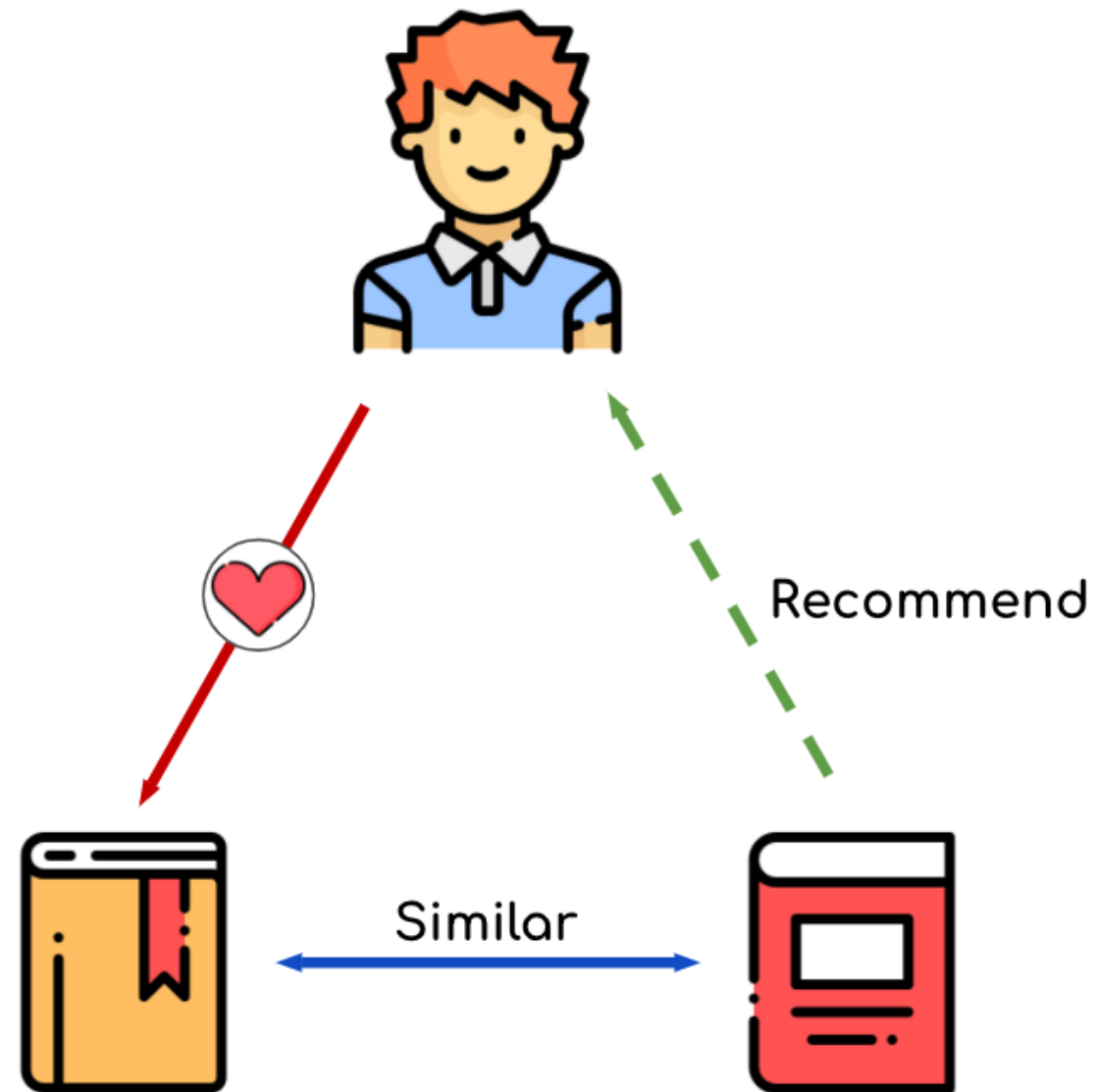


Total price: **$209.02**

[ Add all three to Cart ]

[ Add all three to List ]

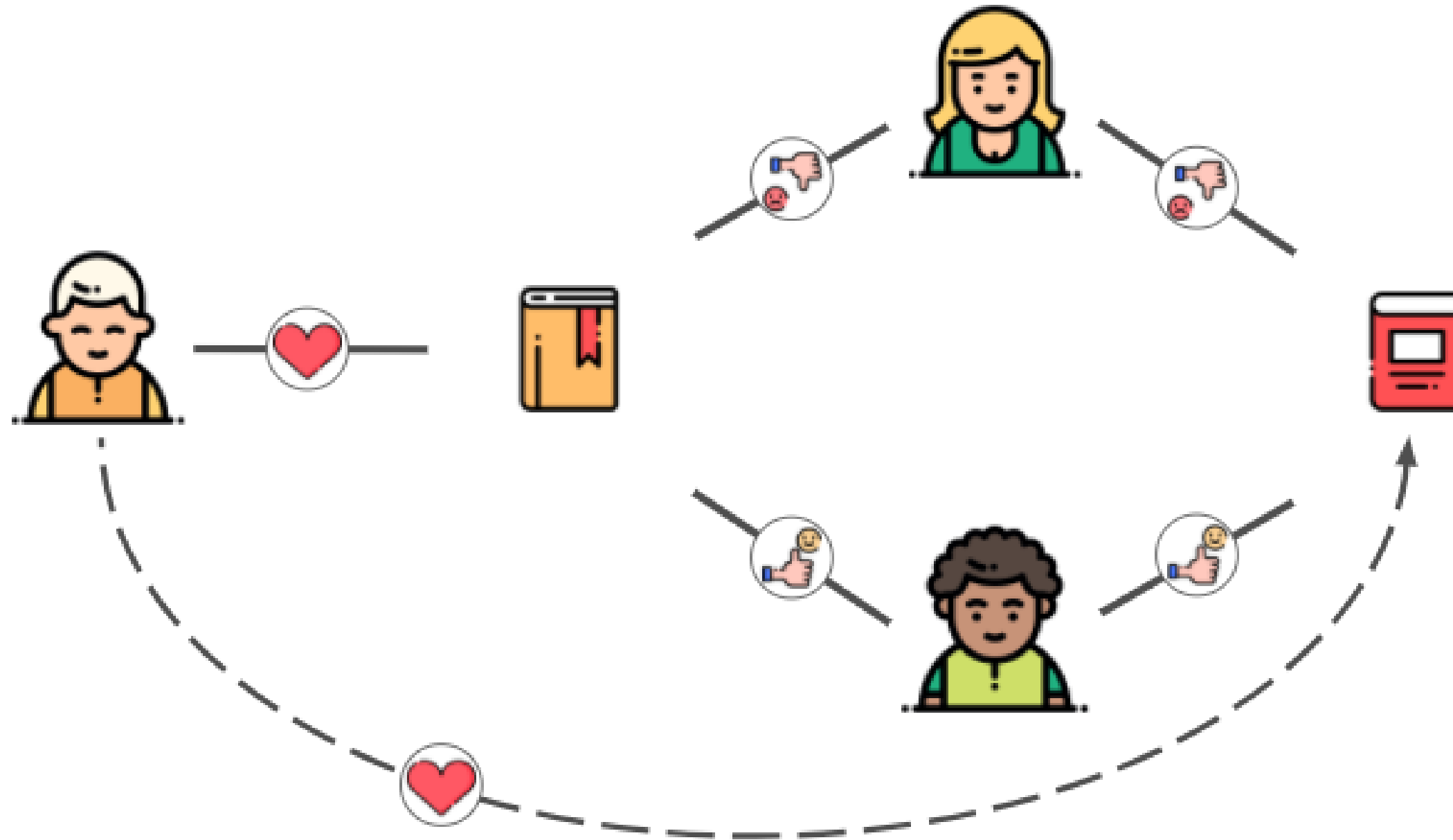ℹ These items are shipped from and sold by different sellers. Show details

☑ **This item:** The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition... by Trevor Hastie Hardcover $71.84

☑ An Introduction to Statistical Learning: with Applications in R (Springer Texts in Statistics) by Gareth James Hardcover $65.39

☑ Pattern Recognition and Machine Learning (Information Science and Statistics) by Christopher M. Bishop Hardcover $71.79

# Content-based models
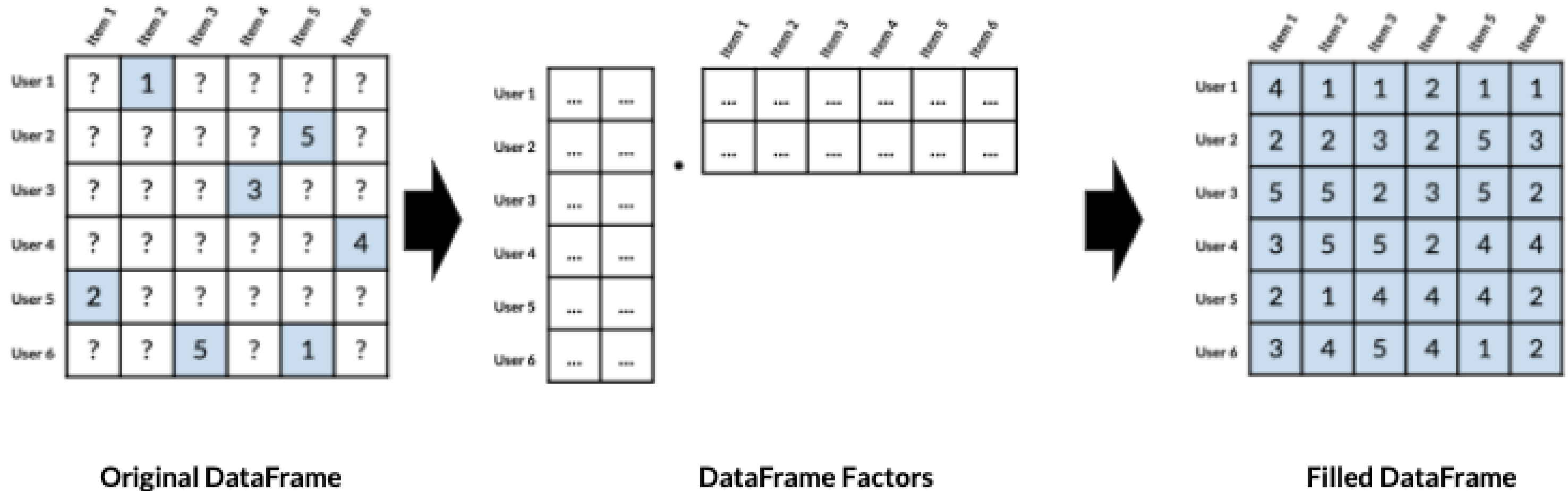
# Collaborative filtering

# Matrix factorization



Original DataFrame

DataFrame Factors

Filled DataFrame

# Congratulations!

## BUILDING RECOMMENDATION ENGINES IN PYTHON