



**Министерство образования и науки Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет
имени Н.Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н.Э. Баумана)**

ФАКУЛЬТЕТ

ИНФОРМАТИКА И СИСТЕМЫ УПРАВЛЕНИЯ

КАФЕДРА

СИСТЕМЫ ОБРАБОТКИ ИНФОРМАЦИИ И УПРАВЛЕНИЯ

**Отчет по лабораторной работе № 4
«Создание рекомендательной модели»
по курсу “Методы машинного обучения”**

**Исполнитель:
Студент группы ИУ5-22М
Желанкина А.С.
22.04.2021**

Москва, 2021

Задание лабораторной работы

Выбрать произвольный набор данных (датасет), предназначенный для построения рекомендательных моделей.

Опираясь на материалы лекции, сформировать рекомендации для одного пользователя (объекта) двумя произвольными способами.

Сравнить полученные рекомендации (если это возможно, то с применением метрик).

Описание датасета

Буккроссинг - это процесс выпуска книг «в дикую природу» для постороннего человека или посредством «контролируемого выпуска» другому участнику BookCrossing и отслеживание того, куда они попадают, с помощью журнальных записей со всего мира.

Экранные формы с текстом программы и примерами её выполнения

```
[ ] import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
import os
from tqdm import tqdm
from gensim.models import Word2Vec
import random
from surprise.accuracy import rmse
from sklearn.metrics.pairwise import cosine_similarity, euclidean_distances, manhattan_distances
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from google.colab import drive
warnings.filterwarnings("ignore")
%matplotlib inline
sns.set(style="ticks")
```

```
[ ] drive.mount("/content/gdrive", force_remount=True)
```

Mounted at /content/gdrive

```
[ ] users = pd.read_csv('/content/gdrive/My Drive/BX-Users.csv',
                        error_bad_lines=False, delimiter=';',
                        encoding = 'ISO-8859-1')
```

```
[ ] books = pd.read_csv('/content/gdrive/My Drive/BX-Books.csv',
                        error_bad_lines=False, delimiter=';',
                        encoding = 'ISO-8859-1')
```

```
b'Skipping line 6452: expected 8 fields, saw 9\nSkipping line 43667: expected 8 fields, saw 10\nSkipping line 51751: expected 8 fields, saw 9\n'
b'Skipping line 92038: expected 8 fields, saw 9\nSkipping line 104319: expected 8 fields, saw 9\nSkipping line 121768: expected 8 fields, saw 9\n'
b'Skipping line 144058: expected 8 fields, saw 9\nSkipping line 150789: expected 8 fields, saw 9\nSkipping line 157128: expected 8 fields, saw 9\nSkipi
b'Skipping line 209388: expected 8 fields, saw 9\nSkipping line 220626: expected 8 fields, saw 9\nSkipping line 227933: expected 8 fields, saw 11\nSki
```

```
[ ] ratings = pd.read_csv('/content/gdrive/My Drive/BX-Book-Ratings.csv',
                          error_bad_lines=False, delimiter=';',
                          encoding = 'ISO-8859-1')
```

```
[ ] print("USERS: ", users.shape, users.columns)
print("BOOKS: ", books.shape, books.columns)
print("RATINGS: ", ratings.shape, ratings.columns)
```

```

USERS: (278858, 3) Index(['User-ID', 'Location', 'Age'], dtype='object')
BOOKS: (271360, 8) Index(['ISBN', 'Book-Title', 'Book-Author', 'Year-Of-Publication', 'Publisher',
                          'Image-URL-S', 'Image-URL-M', 'Image-URL-L'],
                          dtype='object')
RATINGS: (1149780, 3) Index(['User-ID', 'ISBN', 'Book-Rating'], dtype='object')

```

```

[ ] data = pd.merge(ratings, users, on='User-ID', how='inner')
data = pd.merge(data, books, on='ISBN', how='inner')
print("ALL DATA", data.shape, data.columns)

```

```

ALL DATA (1031136, 12) Index(['User-ID', 'ISBN', 'Book-Rating', 'Location', 'Age', 'Book-Title',
                              'Book-Author', 'Year-Of-Publication', 'Publisher', 'Image-URL-S',
                              'Image-URL-M', 'Image-URL-L'],
                              dtype='object')

```

```

[ ] data.head()

```

	User-ID	ISBN	Book-Rating	Location	Age	Book-Title	Book-Author	Year-Of-Publication	Publisher	Image-URL-S	Image-URL-M	Image-URL-L
0	276725	034545104X	0	tyler, texas, usa	NaN	Flesh Tones: A Novel	M. J. Rose	2002	Ballantine Books	http://images.amazon.com/images/P/034545104X.0...	http://images.amazon.com/images/P/034545104X.0...	http://images.amazon.com/images/P/034545104X.0...
1	2313	034545104X	5	cincinnati, ohio, usa	23.0	Flesh Tones: A Novel	M. J. Rose	2002	Ballantine Books	http://images.amazon.com/images/P/034545104X.0...	http://images.amazon.com/images/P/034545104X.0...	http://images.amazon.com/images/P/034545104X.0...
2	6543	034545104X	0	strafford, missouri, usa	34.0	Flesh Tones: A Novel	M. J. Rose	2002	Ballantine Books	http://images.amazon.com/images/P/034545104X.0...	http://images.amazon.com/images/P/034545104X.0...	http://images.amazon.com/images/P/034545104X.0...
3	8680	034545104X	5	st. charles county, missouri, usa	2.0	Flesh Tones: A Novel	M. J. Rose	2002	Ballantine Books	http://images.amazon.com/images/P/034545104X.0...	http://images.amazon.com/images/P/034545104X.0...	http://images.amazon.com/images/P/034545104X.0...
4	10314	034545104X	9	beaverton, oregon, usa	NaN	Flesh Tones: A Novel	M. J. Rose	2002	Ballantine Books	http://images.amazon.com/images/P/034545104X.0...	http://images.amazon.com/images/P/034545104X.0...	http://images.amazon.com/images/P/034545104X.0...

```

[ ] data.info()

```

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 1031136 entries, 0 to 1031135
Data columns (total 12 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   User-ID                             1031136 non-null  int64
1   ISBN                               1031136 non-null  object
2   Book-Rating                         1031136 non-null  int64
3   Location                            1031136 non-null  object
4   Age                                 753301 non-null   float64
5   Book-Title                          1031136 non-null  object
6   Book-Author                         1031135 non-null  object
7   Year-Of-Publication                 1031136 non-null  object
8   Publisher                           1031134 non-null  object
9   Image-URL-S                        1031136 non-null  object
10  Image-URL-M                        1031136 non-null  object
11  Image-URL-L                        1031132 non-null  object
dtypes: float64(1), int64(2), object(9)
memory usage: 102.3+ MB

```

```

[ ] print('Number of books: ', data['ISBN'].nunique())
print('Number of users: ', data['User-ID'].nunique())

```

```
Number of books: 270151
Number of users: 92106
```

```
[ ] median = data["Age"].median()
std = data["Age"].std()
is_null = data["Age"].isnull().sum()
rand_age = np.random.randint(median - std, median + std, size = is_null)
age_slice = data["Age"].copy()
age_slice[pd.isnull(age_slice)] = rand_age
data["Age"] = age_slice
data["Age"] = data["Age"].astype(int)
```

```
[ ] data['Book-Rating'] = data['Book-Rating'].replace(0, None)
```

```
[ ] data[['Book-Author', 'Publisher']] = data[['Book-Author', 'Publisher']].fillna('Unknown')
```

```
[ ] data = data.dropna(axis=0, how='any')
data = data.drop(['Image-URL-S', 'Image-URL-M', 'Image-URL-L'], axis=1)
```

```
[ ] data.isnull().sum()
```

```
User-ID          0
ISBN             0
Book-Rating      0
Location         0
Age              0
Book-Title       0
Book-Author      0
Year-Of-Publication 0
Publisher        0
dtype: int64
```

```
[ ] data['Country'] = data['Location'].apply(lambda row: str(row).split(',')[ -1])
data = data.drop('Location', axis=1)
data['Country'].head()
```

```
0    usa
1    usa
2    usa
3    usa
4    usa
Name: Country, dtype: object
```

```
[ ] data['Year-Of-Publication'] = pd.to_numeric(data['Year-Of-Publication'])
```

```
[ ] df = data
```

```
[ ] df = df[df['Book-Rating'] >= 6]
df.groupby('ISBN')['User-ID'].count().describe()
```

```
count      228988.000000
mean         3.728409
std         12.416574
min          1.000000
25%          1.000000
50%          1.000000
75%          3.000000
max        1206.000000
Name: User-ID, dtype: float64
```

```
[ ] title = df['Book-Title'].values
author = df['Book-Author'].values
publisher = df['Publisher'].values
ISBN = df['ISBN'].values
rating = df['Book-Rating'].values
```

```
[ ] %%time
tfidf = TfidfVectorizer()
title_matrix = tfidf.fit_transform(title)
author_matrix = tfidf.fit_transform(author)
publisher_matrix = tfidf.fit_transform(publisher)
ISBN_matrix = tfidf.fit_transform(ISBN)
```

```
CPU times: user 17.3 s, sys: 167 ms, total: 17.5 s
Wall time: 17.5 s
```

Фильтрация на основе содержания

```
[ ] class SimpleKNNRecommender:

    def __init__(self, X_matrix, X_ids, X_title, X_overview):
        """
        Входные параметры:
        X_matrix - обучающая выборка (матрица объект-признак)
        X_ids - массив идентификаторов объектов
        X_title - массив названий объектов
        X_overview - массив описаний объектов
        """
        #Сохраняем параметры в переменных объекта
        self._X_matrix = X_matrix
        self.df = pd.DataFrame(
            {'id': pd.Series(X_ids, dtype='int'),
             'title': pd.Series(X_title, dtype='str'),
             'overview': pd.Series(X_overview, dtype='str'),
             'dist': pd.Series([], dtype='float')})

    def recommend_for_single_object(self, K: int, \
                                     X_matrix_object, cos_flag = True, manh_flag = False):
        """
        Метод формирования рекомендаций для одного объекта.
        Входные параметры:
        K - количество рекомендуемых соседей
        X_matrix_object - строка матрицы объект-признак, соответствующая объекту
        cos_flag - флаг вычисления косинусного расстояния
        manh_flag - флаг вычисления манхэттэнского расстояния
        Возвращаемое значение: K найденных соседей
        """

        scale = 1000000
        # Вычисляем косинусную близость
        if cos_flag:
            dist = cosine_similarity(self._X_matrix, X_matrix_object)
            self.df['dist'] = dist * scale
            res = self.df.sort_values(by='dist', ascending=False)
            # Не учитываем рекомендации с единичным расстоянием,
            # так как это искомый объект
            res = res[res['dist'] < scale]

        else:
            if manh_flag:
                dist = manhattan_distances(self._X_matrix, X_matrix_object)
            else:
                dist = euclidean_distances(self._X_matrix, X_matrix_object)
            self.df['dist'] = dist * scale
            res = self.df.sort values(by='dist', ascending=True)
```

```
[ ]      # Не учитываем рекомендации с единичным расстоянием,
          # так как это искомый объект
          res = res[res['dist'] > 0.0]

          # Оставляем K первых рекомендаций
          res = res.head(K)
          return res
```

```
[ ]  author_ind = 54
      author[author_ind]
```

'Nicholas Sparks'

```
[ ]  sparks_matrix = author_matrix[author_ind]
```

```
[ ]  skr1 = SimpleKNNRecommender(author_matrix, rating, title, author)
      rec1 = skr1.recommend_for_single_object(15, sparks_matrix)
      rec1
```

		id	title	overview	dist
	442981	10	A Lifelong Passion: Nicholas and Alexandra	Nicholas	679127.939310
[]	442978	10	A Lifelong Passion: Nicholas and Alexandra	Nicholas	679127.939310
	442977	10	A Lifelong Passion: Nicholas and Alexandra	Nicholas	679127.939310
	328753	6	The Discovery of Animal Behaviour	John Sparks	617606.995898
	623204	9	The Last of the Cockleshell Heroes: A World Wa...	William Sparks	570039.764169
	822969	7	The Next Archaeology Workbook	Nicholas David	526087.742446
	516949	9	Cook: The Extraordinary Voyages of Captain Jam...	Nicholas Thomas	496039.823116
	589733	10	The Elephant Man	Christine Sparks	485967.020337
	852180	7	Elephant Man	Christine Sparks	485967.020337
	589732	10	The Elephant Man	Christine Sparks	485967.020337
	281592	7	Veronica	Nicholas Christopher	475879.754871
	390227	9	Veronica: A Novel	Nicholas Christopher	475879.754871

```
[ ]  rec2 = skr1.recommend_for_single_object(15, sparks_matrix, cos_flag = False)
      rec2
```



	id	title	overview	dist
442981	10	A Lifelong Passion: Nicholas and Alexandra	Nicholas	801089.334207
442978	10	A Lifelong Passion: Nicholas and Alexandra	Nicholas	801089.334207
442976	10	A Lifelong Passion: Nicholas and Alexandra	Nicholas	801089.334207
442977	10	A Lifelong Passion: Nicholas and Alexandra	Nicholas	801089.334207
442979	10	A Lifelong Passion: Nicholas and Alexandra	Nicholas	801089.334207
442980	10	A Lifelong Passion: Nicholas and Alexandra	Nicholas	801089.334207
328753	6	The Discovery of Animal Behaviour	John Sparks	874520.444703
623204	9	The Last of the Cockleshell Heroes: A World Wa...	William Sparks	927318.969752
822969	7	The Next Archaeology Workbook	Nicholas David	973562.794640
456102	8	The Sensuous Woman	J	1000000.000000
748786	8	The White Road Westward	B.B.	1000000.000000
748787	8	The Autumn Road to the Isles	B.B.	1000000.000000
781072	8	Cuentos de Encuentros y Desencuentros Amorosos	V. V. A. A.	1000000.000000

Метод на основе сингулярного разложения

```
[ ] def create_utility_matrix(data):
    itemField = 'Book-Title'
    userField = 'User-ID'
    valueField = 'Book-Rating'

    userList = data[userField].tolist()
    itemList = data[itemField].tolist()
    valueList = data[valueField].tolist()

    users = list(set(userList))
    items = list(set(itemList))

    users_index = {users[i]: i for i in range(len(users))}
    pd_dict = {item: [0.0 for i in range(len(users))] for item in items}

    for i in range(0, data.shape[0]):
        item = itemList[i]
        user = userList[i]
        value = valueList[i]
```



```

        pd_dict[item][users_index[user]] = value

    X = pd.DataFrame(pd_dict)
    X.index = users

    itemcols = list(X.columns)
    items_index = {itemcols[i]: i for i in range(len(itemcols))}

    return X, users_index, items_index

```

```
[ ] mini_df = df[0:500]
```

```
[ ] %%time
    user_item_matrix, users_index, items_index = create_utility_matrix(mini_df)
```

```

CPU times: user 3.44 ms, sys: 0 ns, total: 3.44 ms
Wall time: 3.16 ms

```

```
[ ] user_item_matrix
```

```
[ ] user_item_matrix__test = user_item_matrix.iloc[469]
    user_item_matrix__train = user_item_matrix.iloc[:470]
```

```
[ ] %%time
    U, S, VT = np.linalg.svd(user_item_matrix__train.T)
    V = VT.T
```

```

CPU times: user 8.76 ms, sys: 7.17 ms, total: 15.9 ms
Wall time: 28.9 ms

```

```
[ ] Sigma = np.diag(S)
```

```
[ ] r=3
    Ur = U[:, :r]
    Sr = Sigma[:r, :r]
    Vr = V[:, :r]
```

```
[ ] test_user = np.mat(user_item_matrix__test.values)
    test_user.shape, test_user
```

```
((1, 7), matrix([[8., 0., 0., 0., 0., 0., 0.])))
```

```
[ ] tmp = test_user * Ur * np.linalg.inv(Sr)
    tmp
```

```
matrix([[ -0.013692 ,  0.07842731,  0.02007541]])
```

```
[ ] test_user_result = np.array([tmp[0,0], tmp[0,1], tmp[0,2]])
    test_user_result
```

```
array([ -0.013692 ,  0.07842731,  0.02007541])
```

```
[ ] # Вычисляем косинусную близость между текущим пользователем
    # и остальными пользователями
    cos_sim = cosine_similarity(Vr, test_user_result.reshape(1, -1))
    cos_sim[:10]
```

```
array([[ -0.05145611],
       [ -0.02704594],
       [  1.          ],
       [ -0.02704594],
       [  0.87936934],
       [ -0.02704594],
```

```
[ ] # Преобразуем размерность массива
    cos_sim_list = cos_sim.reshape(-1, cos_sim.shape[0])[0]
    cos_sim_list[:10]
```

```
array([ -0.05145611, -0.02704594,  1.          , -0.02704594,  0.87936934,
        -0.02704594,  1.          , -0.02704594,  1.          ,  0.87936934])
```

```
[ ] # Находим наиболее близкого пользователя
    recommended_user_id = np.argsort(-cos_sim_list)[0]
    recommended_user_id
```

254

```
[ ] # Получение названия фильма
    userId_list = list(user_item_matrix.columns)
    def book_name_by_user(ind):
        try:
            userId = userId_list[ind]
            flt_links = mini_df[mini_df['User-ID'] == userId]
            #tmdbId = int(flt_links['tmdbId'].values[0])
            #md_links = df_md[df_md['id'] == tmdbId]
            res = mini_df['Book-Title'].values[0]
            return res
```

```
except:  
    return ''
```



```
i=1  
for idx, item in enumerate(np.ndarray.flatten(np.array(test_user))):  
    if item > 0:  
        book_title = book_name_by_user(idx)  
        print('{} - {} - {}'.format(idx, book_title, item))
```

0 - Flesh Tones: A Novel - 8.0