Рекомендательные системы в продажах # В интернет-магазине есть данные по пользовательским сессиям, каждая из которых содержит список # просмотренных товаров (id) и список купленных товаров. # Необходимо на основе этих данных создать рекомендательную систему, предлагающую пользователю # k - товаров на основе его просмотров import pandas as pd In [1]: import numpy as np with open('Data/Coursera/coursera\_sessions\_train.txt', 'r') as f: sess\_train = f.read().splitlines() with open('Data/Coursera/coursera\_sessions\_test.txt', 'r') as f: sess\_test = f.read().splitlines() # представленные данные: каждая строка - список просмотренных товаров (id) пользователем за одну сессию, а через ";" - купленный sess\_train[8:12] Out[5]: ['71,72,73,74;', '76,77,78;', '84,85,86,87,88,89,84,90,91,92,93,86;86', '114,77,115,116,117,118,119,120,121,120,122,123,124;'] # Проект построен следующим образом # 1. Создаём матрицы частот просматриваемых и покупаемых товаров # (товар - кол-во просмотров и товар - кол-во покупок), сортируем м-цы по убыванию частот. # Т.о., получили м-цу популярности просмотров и м-цу популярности покупок # 2. Для текущей сессии просмотров [ $id_1$ , $id_2$ ,..., $id_n$ ] делаем рекомендации # на k товаров из этого списка (k <= n), стоящие в m-це популярности ПРОСМОТРОВ выше оставшихся n-k просмотров # 3. Для текущей сессии просмотров  $[id_1,id_2,...,id_n]$  делаем рекомендации # на k товаров из этого списка (k <= n), стоящие в m-це популярности ПОКУПОК выше оставшихся n - k просмотров # 4. - 6. Проверяем все то же самое на тестовой выборке (используя, разумеется, м-цы популярности тренировочной выборки) # В итоге, значения метрик для рекомендаций по покупкам ведут себя на тестовой выборке немного лучше, # при этом значения метрик для тестовой выборке очень близки к значениям метрик по тренировочной выборке # для рекомендаций по просмотрам. 1. На обучении постройте частоты появления id в просмотренных и в купленных In [ ]: #(id может несколько раз появляться в просмотренных, все появления надо учитывать) # разделяем м-цу сессий на просмотренные id и купленные id In [3]: # для обучающей выборки sess\_train\_lp = [] for sess in sess\_train: look\_items, pur\_items = sess.split(';') # применяем ф-цию int() ко всем id сессии look\_items = list(map(int, look\_items.split(','))) if len(pur\_items) > 0: pur\_items = list(map(int, pur\_items.split(','))) else: pur\_items = [] sess\_train\_lp.append([look\_items, pur\_items]) # для тестовой выборки sess\_test\_lp = [] for sess in sess\_test: look\_items, pur\_items = sess.split(';') look\_items = list(map(int, look\_items.split(','))) if len(pur\_items) > 0: pur\_items = list(map(int, pur\_items.split(','))) else: pur\_items = [] sess\_test\_lp.append([look\_items, pur\_items]) sess\_train\_lp[8:12] In [10]: Out[10]: [[[71, 72, 73, 74], []], [[76, 77, 78], []], [[84, 85, 86, 87, 88, 89, 84, 90, 91, 92, 93, 86], [86]], [[114, 77, 115, 116, 117, 118, 119, 120, 121, 120, 122, 123, 124], []]] # ДЛЯ ПРОСМОТРЕННЫХ In [4]: # сначала все просмотренные id загоняем в список sess\_train\_lid = [] for i in range(len(sess\_train\_lp)): for id in sess\_train\_lp[i][0]: sess\_train\_lid.append(id) # теперь столбец уникальных id со счетчиком sess\_train\_l\_cnt = np.transpose(np.unique(sess\_train\_lid, return\_counts=True)) # структура столбца: i-й элемент является двумерной строкой [id, кол-во раз которое он встречается] In [22]: sess\_train\_l\_cnt Out[22]: array([[ 6], 9], [102804, [102805, 1]], dtype=int64) [102806, # ДЛЯ КУПЛЕННЫХ # сначала все просмотренные id загоняем в список sess\_train\_pid = [] for i in range(len(sess\_train\_lp)): for id in sess\_train\_lp[i][1]: sess\_train\_pid.append(id) # теперь столбец уникальных id со счетчиком sess\_train\_p\_cnt = np.transpose(np.unique(sess\_train\_pid, return\_counts=True)) sess\_train\_p\_cnt In [24]: Out[24]: array([[ 1], 1], [102417, [102462, [102646, 1]], dtype=int64) # Сортируем полученные столбцы по счетчику (запись по убыванию) # sess\_train\_l\_cnt[:,1].argsort() - список индексов строк нашего двумерного массива по возрастанию значения счетчика # sess\_train\_l\_cnt[sess\_train\_l\_cnt[:,1].argsort()] - отсортированный в порядке возрастания счетчика массив # onepaция среза - list[<start>:<stop>:<step>], [::-1] - просто вернет список с конца до начала sess\_train\_l\_cnt = sess\_train\_l\_cnt[sess\_train\_l\_cnt[:,1].argsort()][::-1] sess\_train\_p\_cnt = sess\_train\_p\_cnt[sess\_train\_p\_cnt[:,1].argsort()][::-1] In [40]: # итоговая м-ца купленных товаров: id товара - сколько раз он был куплен sess\_train\_p\_cnt Out[40]: array([[ 158, 14], 12], 204, 11], 73, [38189, 1], [38177, 1], 1]], dtype=int64) [ 5, 2. Алгоритм рекомендаций - сортировка просмотренных id по популярности (частота появления в просмотренных) # ф-ция расчета метрик 'Precision@k' и 'Recall@k' для k-рекомендаций In [7]: # 'reccomendations' - список 'k' рекомендуемых товаров для текущей сессии, в которой был совершен список покупок 'session' def prec\_rec\_metrics(session, reccomendations, k): purchase = 0 for ind in reccomendations: if ind in session: purchase += 1 precision = purchase / k recall = purchase / len(session) return(precision, recall) # столбцы купленных и просмотренных товаров sess\_train\_p = [row[1] for row in sess\_train\_lp] sess\_train\_l = [row[0] for row in sess\_train\_lp] %%time In [9]: # Считаем метрики 'Precision@k' и 'Recall@k' на обучающей выборке для k=1 и k=5prec\_at\_1\_tr\_l, rec\_at\_1\_tr\_l = [], [] prec\_at\_5\_tr\_l, rec\_at\_5\_tr\_l = [], [] k1, k5 = 1, 5for i, sess\_p in enumerate(sess\_train\_p): # не обрабатываем сессии без покупок if sess\_p == []: continue # просмотренные ід для сессии і sess\_l = sess\_train\_l[i]  $l_{ind_sess} = []$ for j in range(len(sess\_1)): # выбираем индекс (!) столбца-счетчика встречаемости просмотренных ід, # элемент в котором соответствует значению текущего id сессии. l\_ind\_sess.append(np.where(sess\_train\_l\_cnt[:,0] == sess\_l[j])[0][0]) # список индексов м-цы-счетчика купленных товаров для купленных товаров из текущей сессии l\_ind\_sess\_sorted = np.unique(l\_ind\_sess) # k1 рекомендаций num\_of\_recs\_k1 = min(k1, len(sess\_1)) if num\_of\_recs\_k1 == 0: continue # берем первые Min(k1, кол-во id в текущей сессии) кол-во id из текущей сессии,# согласно сортировке по частоте просмотров по всей обучающей выборке recs\_k1 = sess\_train\_l\_cnt[l\_ind\_sess\_sorted[:num\_of\_recs\_k1],0] # расчет k1-метрик prec\_1, rec\_1 = prec\_rec\_metrics(sess\_p, recs\_k1, k1) prec\_at\_1\_tr\_l.append(prec\_1) rec\_at\_1\_tr\_l.append(rec\_1) # k5 рекомендаций num\_of\_recs\_k5 = min(k5, len(sess\_1)) if num of recs k5 == 0: continue recs\_k5 = sess\_train\_l\_cnt[l\_ind\_sess\_sorted[:num\_of\_recs\_k5],0] # расчет k5-метрик prec\_5, rec\_5 = prec\_rec\_metrics(sess\_p, recs\_k5, k5) prec\_at\_5\_tr\_l.append(prec\_5) rec\_at\_5\_tr\_l.append(rec\_5) Wall time: 8.23 s avg\_prec\_at\_1\_tr\_l = np.mean(prec\_at\_1\_tr\_l) In [10]: avg\_rec\_at\_1\_tr\_l = np.mean(rec\_at\_1\_tr\_l) avg\_prec\_at\_5\_tr\_1 = np.mean(prec\_at\_5\_tr\_1) avg\_rec\_at\_5\_tr\_l = np.mean(rec\_at\_5\_tr\_l) r1 = round(avg\_rec\_at\_1\_tr\_1, 2) In [13]: p1 = round(avg\_prec\_at\_1\_tr\_1, 2) r5 = round(avg\_rec\_at\_5\_tr\_1, 2) p5 = round(avg\_prec\_at\_5\_tr\_1, 2) In [24]: # будем заносить все результаты этого проекта в единый датафрейм 'metrics' metrics = pd.DataFrame({'precision':p1, 'recall': r1},index={'Train\_look\_k=1'}) metrics.loc['Train\_look\_k=5'] = [p5,r5] 3. Алгоритм рекомендаций - сортировка просмотренных id по покупаемости (частота появления в покупках) In [25]: %%time # Считаем метрики 'Precision@k' и 'Recall@k' на обучающей выборке для k=1 и k=5prec\_at\_1\_tr\_p, rec\_at\_1\_tr\_p = [], [] prec\_at\_5\_tr\_p, rec\_at\_5\_tr\_p = [], [] k1, k5 = 1, 5for i, sess\_p in enumerate(sess\_train\_p): # не обрабатываем сессии без покупок if sess\_p == []: continue # просмотренные id для сессии i sess\_l = sess\_train\_l[i] l\_ind\_sess = [] for j in range(len(sess\_l)): # выбираем индекс столбца-счетчика встречаемости просмотренных id, # элемент в котором соответствует значению текущего id сессии # здесь первый [0] выбирает первую строку в найденном однострочном элементе, # второй [0] выбирает первый элемент строки, являющимся как раз индексом if sess\_l[j] not in sess\_train\_p\_cnt[:,0]: continue l\_ind\_sess.append(np.where(sess\_train\_p\_cnt[:,0] == sess\_l[j])[0][0]) l\_ind\_sess\_sorted = np.unique(l\_ind\_sess) # к1 рекомендаций num\_of\_recs\_k1 = min(k1, len(sess\_1), len(l\_ind\_sess\_sorted)) if num\_of\_recs\_k1 == 0: continue # берем первые Min(k1, кол-во id в текущей сессии) кол-во id из текущей сессии,# согласно сортировке по частоте просмотров по всей обучающей выборке recs\_k1 = sess\_train\_p\_cnt[l\_ind\_sess\_sorted[:num\_of\_recs\_k1],0] # расчет k1-метрик prec\_1, rec\_1 = prec\_rec\_metrics(sess\_p, recs\_k1, k1) prec\_at\_1\_tr\_p.append(prec\_1) rec\_at\_1\_tr\_p.append(rec\_1) # k5 рекомендаций num\_of\_recs\_k5 = min(k5, len(sess\_1), len(l\_ind\_sess\_sorted)) if num\_of\_recs\_k5 == 0: continue recs\_k5 = sess\_train\_p\_cnt[l\_ind\_sess\_sorted[:num\_of\_recs\_k5],0] # расчет k5-метрик prec\_5, rec\_5 = prec\_rec\_metrics(sess\_p, recs\_k5, k5) prec\_at\_5\_tr\_p.append(prec\_5) rec\_at\_5\_tr\_p.append(rec\_5) Wall time: 1.74 s avg\_prec\_at\_1\_tr\_p = np.mean(prec\_at\_1\_tr\_p) In [26]: avg\_rec\_at\_1\_tr\_p = np.mean(rec\_at\_1\_tr\_p) avg\_prec\_at\_5\_tr\_p = np.mean(prec\_at\_5\_tr\_p) avg\_rec\_at\_5\_tr\_p = np.mean(rec\_at\_5\_tr\_p) r1 = round(avg\_rec\_at\_1\_tr\_p, 2) p1 = round(avg\_prec\_at\_1\_tr\_p, 2) r5 = round(avg\_rec\_at\_5\_tr\_p, 2) p5 = round(avg\_prec\_at\_5\_tr\_p, 2) metrics.loc['Train\_purch\_k=1'] = [p1,r1] In [27]: metrics.loc['Train\_purch\_k=5'] = [p5,r5] 4. На тестовой выборке построим частоты появления id в просмотренных и в купленных # Формируем столбец просмотренных товаров за все тренировочные сессии In [53]: sess\_test\_l = [row[0] for row in sess\_test\_lp] sess\_test\_l\_np = [] for sess in sess\_test\_1: for idd in sess: sess\_test\_l\_np.append(idd) sess\_test\_l\_np = np.array(sess\_test\_l\_np) # Формируем столбец купленных товаров за все тренировочные сессии sess\_test\_p = [row[1] for row in sess\_test\_lp] sess\_test\_p\_np = [] for sess in sess\_test\_p: for idd in sess: sess\_test\_p\_np.append(idd) sess\_test\_p\_np = np.array(sess\_test\_p\_np) 5. Алгоритм рекомендаций на тестовой выборке - сортировка просмотренных ід по популярности (частота появления в просмотренных) %%time In [57]: # метрики по тестовым данным prec\_at\_1\_tst\_l, rec\_at\_1\_tst\_l = [], [] prec\_at\_5\_tst\_l, rec\_at\_5\_tst\_l = [], [] k1, k5 = 1, 5for i, sess\_p in enumerate(sess\_test\_p): # пропускаем сессии без покупок if sess\_p == []: continue # проходимся по id из просмотренных товаров sess\_l = sess\_test\_l[i] # выбираем индекс столбца-счетчика встречаемости просмотренных ід в тренировочной приоритетной матрице l\_ind\_sess = []  $new_ids = []$ for j in range(len(sess\_1)): if sess\_l[j] not in sess\_train\_l\_cnt[:,0]: new\_ids.append(sess\_l[j]) continue l\_ind\_sess.append(np.where(sess\_train\_l\_cnt[:,0] == sess\_l[j])[0][0]) l\_ind\_sess\_sorted = np.unique(l\_ind\_sess) # к1 рекомендации num\_of\_recs\_k1 = min(k1, len(sess\_1)) if num\_of\_recs\_k1 == 0: continue # добавляем это условие, поскольку все просмотренные в сессии товары могут не быть в тренировочной м-це if l\_ind\_sess != []: recs\_k1 = sess\_train\_l\_cnt[l\_ind\_sess\_sorted[:num\_of\_recs\_k1],0] else:  $recs_k1 = []$ # здесь объединяем отсортированные 'в порядке убывания приоритета приоритетной м-цы тренировочных данных' # рекомендуемые товары и новые товары, не присутствующие ранее в приоритетной м-це тренировочных данных recs\_k1 = np.concatenate((np.array(recs\_k1, dtype='int64'), np.unique(np.array(new\_ids, dtype='int64'))))[:num\_of\_recs\_k1] # k1 метрики prec\_1, rec\_1 = prec\_rec\_metrics(sess\_p, recs\_k1, k1) prec\_at\_1\_tst\_l.append(prec\_1) rec\_at\_1\_tst\_l.append(rec\_1) # k5 рекомендации num\_of\_recs\_k5 = min(k5, len(sess\_1)) if num\_of\_recs\_k5 == 0: continue if l\_ind\_sess != []: recs\_k5 = sess\_train\_l\_cnt[l\_ind\_sess\_sorted[:num\_of\_recs\_k5],0] else:  $recs_k5 = []$ recs\_k5 = np.concatenate((np.array(recs\_k5, dtype='int64'), np.unique(np.array(new\_ids, dtype='int64'))))[:num\_of\_recs\_k5] # k5 метрики prec\_5, rec\_5 = prec\_rec\_metrics(sess\_p, recs\_k5, k5) prec\_at\_5\_tst\_l.append(prec\_5) rec\_at\_5\_tst\_l.append(rec\_5) Wall time: 12.9 s avg\_prec\_at\_1\_tst\_l = np.mean(prec\_at\_1\_tst\_l) In [59]: avg\_rec\_at\_1\_tst\_l = np.mean(rec\_at\_1\_tst\_l) avg\_prec\_at\_5\_tst\_1 = np.mean(prec\_at\_5\_tst\_1) avg\_rec\_at\_5\_tst\_1 = np.mean(rec\_at\_5\_tst\_1) r1 = round(avg\_rec\_at\_1\_tst\_1, 2) p1 = round(avg\_prec\_at\_1\_tst\_l, 2) r5 = round(avg\_rec\_at\_5\_tst\_1, 2) p5 = round(avg\_prec\_at\_5\_tst\_1, 2) metrics.loc['Test\_look\_k=1'] = [p1,r1] In [60]: metrics.loc['Test\_look\_k=5'] = [p5,r5] 6. Алгоритм рекомендаций на тестовой выборке - сортировка просмотренных іd по покупаемости (частота появления в покупках) In [84]: prec\_at\_1\_tst\_p, rec\_at\_1\_tst\_p = [], [] prec\_at\_5\_tst\_p, rec\_at\_5\_tst\_p = [], [] k1, k5 = 1, 5for i, sess\_p in enumerate(sess\_test\_p): if sess\_p == []: continue sess\_l = sess\_test\_l[i] l\_ind\_sess = [] new\_ids = [] for j in range(len(sess\_1)): if sess\_l[j] not in sess\_train\_p\_cnt[:,0]: new\_ids.append(sess\_l[j]) continue l\_ind\_sess.append(np.where(sess\_train\_p\_cnt[:,0] == sess\_l[j])[0][0]) l\_ind\_sess\_sorted = np.unique(l\_ind\_sess) # k1 recommendations num\_of\_recs\_k1 = min(k1, len(sess\_l)) if num\_of\_recs\_k1 == 0: continue if l\_ind\_sess != []: recs\_k1 = sess\_train\_p\_cnt[l\_ind\_sess\_sorted[:num\_of\_recs\_k1],0] else:  $recs_k1 = []$ # химичим тут recs\_k1 = np.concatenate((np.array(recs\_k1, dtype='int64'), np.unique(np.array(new\_ids, dtype='int64'))))[:num\_of\_recs\_k1] # k1 метрики prec\_1, rec\_1 = prec\_rec\_metrics(sess\_p, recs\_k1, k1) prec\_at\_1\_tst\_p.append(prec\_1) rec\_at\_1\_tst\_p.append(rec\_1) # k5 рекомендации num\_of\_recs\_k5 = min(k5, len(sess\_1)) if num\_of\_recs\_k5 == 0: continue if l\_ind\_sess != []: recs\_k5 = sess\_train\_p\_cnt[l\_ind\_sess\_sorted[:num\_of\_recs\_k5],0] else:  $recs_k5 = []$ recs\_k5 = np.concatenate((np.array(recs\_k5, dtype='int64'), np.unique(np.array(new\_ids, dtype='int64'))))[:num\_of\_recs\_k5] # k5 метрики prec\_5, rec\_5 = prec\_rec\_metrics(sess\_p, recs\_k5, k5) prec\_at\_5\_tst\_p.append(prec\_5) rec\_at\_5\_tst\_p.append(rec\_5) avg\_prec\_at\_1\_tst\_p = np.mean(prec\_at\_1\_tst\_p) avg\_rec\_at\_1\_tst\_p = np.mean(rec\_at\_1\_tst\_p) avg\_prec\_at\_5\_tst\_p = np.mean(prec\_at\_5\_tst\_p) avg\_rec\_at\_5\_tst\_p = np.mean(rec\_at\_5\_tst\_p) r1 = round(avg\_rec\_at\_1\_tst\_p, 2) p1 = round(avg\_prec\_at\_1\_tst\_p, 2) r5 = round(avg\_rec\_at\_5\_tst\_p, 2) p5 = round(avg\_prec\_at\_5\_tst\_p, 2) Answer 4: 0.42 0.49 0.80 0.20 metrics.loc['Test\_purch\_k=1'] = [p1,r1] In [86]: metrics.loc['Test\_purch\_k=5'] = [p5,r5] 7. Анализ результатов metrics Out[87]: precision recall Train\_look\_k=1 0.51 0.44 0.21 0.83 Train\_look\_k=5 Train\_purch\_k=1 0.79 0.68 0.25 0.93 Train\_purch\_k=5 0.48 0.42 Test\_look\_k=1 Test\_look\_k=5 0.2 0.8 0.49 0.42 Test\_purch\_k=1 Test\_purch\_k=5 0.20 0.80 metrics.loc[['Train\_look\_k=1','Test\_purch\_k=1']] Out[93]: precision recall Train\_look\_k=1 0.51 0.44 Test\_purch\_k=1 0.49 0.42 metrics.loc[['Train\_look\_k=5','Test\_purch\_k=5']] precision recall Out[94]: Train\_look\_k=5 0.83 0.21 Test\_purch\_k=5 0.20 0.80 metrics.loc[['Train\_look\_k=1','Test\_look\_k=1']] Out[95]: precision recall Train\_look\_k=1 0.51 0.44 Test\_look\_k=1 0.48 0.42 metrics.loc[['Train\_look\_k=5','Test\_look\_k=5']] Out[96]: precision recall Train\_look\_k=5 0.21 0.83 Test\_look\_k=5 0.2 8.0 metrics.loc[['Train\_purch\_k=1','Test\_purch\_k=1']] Out[97]: precision recall Train\_purch\_k=1 0.79 0.68 Test\_purch\_k=1 0.49 0.42 metrics.loc[['Train\_purch\_k=5','Test\_purch\_k=5']] Out[98]: precision recall Train\_purch\_k=5 0.25 0.93 Test\_purch\_k=5 0.20 0.80