

PROGRESS REPORT

PREDICTION OF POTENTIAL DRUG-TARGET INTERACTIONS USING PROBABILISTIC MATRIX FACTORIZATION

INFSCI 2725 DATA ANALYTICS

JULY 21, 2019

GROUP MEMBERS:

PREM RAJGOPAL: PRR47@pitt.edu

TSHERING SHERPA: TS12@pitt.edu

ZHENG HAN: ZHH41@pitt.edu

Problem Identification

We chose to study the unknown interactions between existing drugs and the well-known proteins in human or animals. This prediction will provide a guideline for the future drug discovery research and also help to explain the side effects of certain drugs.

Questions To Be Answered

Q1. For the currently reported interactions, find out the average number of targets for the drugs, and list the top 100 drugs according to the number of targets. Draw a histogram for the top 100 drugs in terms of the number of targets.

Q2. Predict the unknown interactions between drugs and targets. In other words, look for more drug-target pairs. Predict the interaction probability of each newly found drug-target pair.

Q3. The reported interactions serve as the positive samples. A certain number of negative samples need to be generated for training the model. The effect of the number of negative samples on the RMSE and AUC of the model will be studied.

What We Have Accomplished So Far

We have successfully answered Q1. For Q2 and Q3, we generated preliminary results with limited accuracy, and we are still trying to improve the accuracy.

Preliminary Results

Q1:

The original data set downloaded from DrugBank database was imported into a Jupyter Notebook and it looks like below:

```
data = pd.read_csv('uniprot_links_2.csv')
data.head()
```

	DrugBankID	Name	Type	UniProtID	UniProtName
0	DB00001	Lepirudin	BiotechDrug	P00734	Prothrombin
1	DB00002	Cetuximab	BiotechDrug	P00533	Epidermal growth factor receptor
2	DB00002	Cetuximab	BiotechDrug	O75015	Low affinity immunoglobulin gamma Fc region re...
3	DB00002	Cetuximab	BiotechDrug	P00736	Complement C1r subcomponent
4	DB00002	Cetuximab	BiotechDrug	P02745	Complement C1q subcomponent subunit A

By using `data.describe()`, we know that there are 7370 drugs represented by “DrugBankID” and 4763 proteins (targets) represented by “UniProtID” as shown below:

```
data.describe()
```

	DrugBankID	Name	Type	UniProtID	UniProtName
count	20061	20061	20061	20061	20061
unique	7370	7370	2	4763	4255
top	DB12010	Fostamatinib	SmallMoleculeDrug	P24941	Cyclin-dependent kinase 2
freq	306	306	19134	136	136

By grouping “DrugBankID”, we know how many targets for each drug. And then by using `mean()` function, we know the average number of targets for the drugs is about 2.72.

```
grouped=data.groupby(['DrugBankID']).count()
grouped.head()
```

	Name	Type	UniProtID	UniProtName
DrugBankID				
DB00001	1	1	1	1
DB00002	12	12	12	12
DB00004	3	3	3	3
DB00005	14	14	14	14
DB00006	1	1	1	1

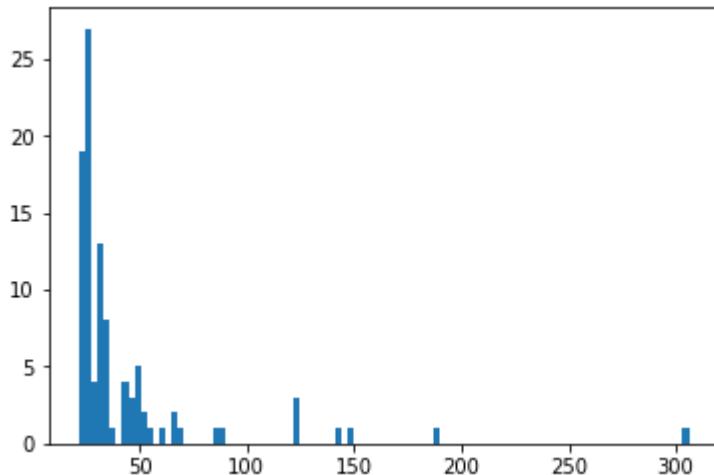
```
grouped['UniProtID'].mean()
```

2.7219810040705563

Then the drugs are sorted according to the number of targets and the top 100 drugs are displayed. A histogram is drawn to show the distribution of the number of targets.

```
first100=grouped.sort_values('Name', ascending=False).head(100)
first100.head()
```

	Name	Type	UniProtID	UniProtName
DrugBankID				
DB12010	306	306	306	306
DB11638	189	189	189	189
DB09130	147	147	147	147
DB00157	144	144	144	144
DB01593	124	124	124	124



Q2:

The Probabilistic Matrix Factorization (PMF) model [1] is used for our prediction. We found the implementation of PMF from Github [2]. This author used PMF to predict the user's rating on the movies that the user has not watched. The rating score ranges from 1 to 5 in his work. The prediction of drug and target interaction is similar to the prediction of user and movie relationship. But our scale is 0 to 1 with 0 indicating no interaction and 1 indicating interaction with a probability of 100%.

In order to use the PMF model, the drug ID, target ID and interaction scale have to be changed to integer. Below is what we did to change the labels.

All the DrugBankIDs are saved in a dictionary with the DrugBankID as the *key* and the integer value starting from 1 as the *value*:

```

dict1 = {}
counter = 1
for item in set(data['DrugBankID']):
    dict1[item] = counter
    counter += 1
print(dict1)

{'DB00913': 1, 'DB05890': 2, 'DB03028': 3, 'DB07312': 4, 'DB00687': 5, 'DB03876': 6, 'DB04891': 7, 'DB14562': 8, 'DB12442': 9, 'DB08049': 10, 'DB04852': 11, 'DB07806': 12, 'DB02767': 13, 'DB03254': 14, 'DB04459': 15, 'DB07440': 16, 'DB04547': 17, 'DB07733': 18, 'DB02120': 19, 'DB14541': 20, 'DB09107': 21, 'DB03961': 22, 'DB03568': 23, 'DB05630': 24, 'DB02954': 25, 'DB07858': 26, 'DB07641': 27, 'DB03590': 28, 'DB00491': 29, 'DB07409': 30, 'DB00357': 31, 'DB03887': 32, 'DB07020': 33, 'DB07696': 34, 'DB00377': 35, 'DB00590': 36, 'DB08096': 37, 'DB03921': 38, 'DB06247': 39, 'DB04743': 40, 'DB09510': 41, 'DB05271': 42, 'DB07111': 43, 'DB07748': 44, 'DB00723': 45, 'DB08145': 46, 'DB08607': 47, 'DB03022': 48, 'DB03827': 49, 'DB06938': 50, 'DB06345': 51, 'DB04284': 52, 'DB04482': 53, 'DB05434': 54, 'DB03766': 55, 'DB01145': 56, 'DB07628': 57, 'DB11580': 58, 'DB02741': 59, 'DB00957': 60, 'DB08637': 61, 'DB08299': 62, 'DB01613': 63, 'DB04223': 64, 'DB08810': 65, 'DB09026': 66, 'DB03210': 67, 'DB06237': 68, 'DB04081': 69, 'DB05967': 70, 'DB00750': 71, 'DB04672': 72, 'DB00407': 73, 'DB04562': 74, 'DB03723': 75, 'DB01980': 76, 'DB11397': 77, 'DB08614': 78, 'DB07955': 79, 'DB00398': 80, 'DB07249': 81, 'DB01419': 82, 'DB01268': 83, 'DB08067': 84, 'DB07276': 85, 'DB08264': 86, 'DB08357': 87, 'DB03240': 88, 'DB00534': 89, 'DB04629': 90, 'DB06691': 91, 'DB13943': 92, 'DB04271': 93, 'DB13165': 94, 'DB04235': 95, 'DB14004': 96, 'DB08727': 97, 'DB07205': 98, 'DB01158': 99, 'DB02790': 100, 'DB08238': 101, 'DB08466': 102, 'DB12872': 103, 'DB083658': 104, 'DB083358': 105, 'DB08062': 106, 'DB03577': 107, 'DB08360': 108, 'DB05448': 109, 'DB03331': 110, 'DB04793': 111, 'DB03556': 112, 'DB06977': 113, 'DB11071': 114, 'DB06710': 115, 'DB05276': 116, 'DB087889': 117, 'DB00485': 118, 'DB02366': 119, 'DB03480': 120, 'DB02494': 121, 'DB087602': 122, 'DB00947': 123, 'DB03821': 124, 'DB00290': 125, 'DB00474': 126, 'DB02425': 127, 'DB03181': 128, 'DB05944': 129, 'DB00613': 130, 'DB07268': 131, 'DB00940': 132, 'DB02375': 133, 'DB07968': 134, 'DB06674': 135, 'DB07175': 136, 'DB07219': 137, 'DB07326': 138, 'DB06226': 139, 'DB07888': 140, 'DB03668': 141, 'DB06263': 142, 'DB04261': 143, 'DB03416': 144, 'DB05023': 145, 'DB07507': 146, 'DB04040': 147, 'DB02955': 148, 'DB088295': 149, 'DB008895': 150, 'DB06446': 151, 'DB04903': 152, 'DB04895': 153, 'DB03868': 154, 'DB00961': 155, 'DB02170': 156, 'DB012141': 157, 'DB022571': 158, 'DB002551': 159, 'DB022411': 160, 'DB066281': 162, 'DB071151': 163, 'DB03822': 164, 'DB04041': 165, 'DB04042': 166, 'DB04043': 167, 'DB04044': 168, 'DB04045': 169, 'DB04046': 170, 'DB04047': 171, 'DB04048': 172, 'DB04049': 173, 'DB04050': 174, 'DB04051': 175, 'DB04052': 176, 'DB04053': 177, 'DB04054': 178, 'DB04055': 179, 'DB04056': 180, 'DB04057': 181, 'DB04058': 182, 'DB04059': 183, 'DB04060': 184, 'DB04061': 185, 'DB04062': 186, 'DB04063': 187, 'DB04064': 188, 'DB04065': 189, 'DB04066': 190, 'DB04067': 191, 'DB04068': 192, 'DB04069': 193, 'DB04070': 194, 'DB04071': 195, 'DB04072': 196, 'DB04073': 197, 'DB04074': 198, 'DB04075': 199, 'DB04076': 200, 'DB04077': 201, 'DB04078': 202, 'DB04079': 203, 'DB04080': 204, 'DB04081': 205, 'DB04082': 206, 'DB04083': 207, 'DB04084': 208, 'DB04085': 209, 'DB04086': 210, 'DB04087': 211, 'DB04088': 212, 'DB04089': 213, 'DB04090': 214, 'DB04091': 215, 'DB04092': 216, 'DB04093': 217, 'DB04094': 218, 'DB04095': 219, 'DB04096': 220, 'DB04097': 221, 'DB04098': 222, 'DB04099': 223, 'DB040100': 224, 'DB040101': 225, 'DB040102': 226, 'DB040103': 227, 'DB040104': 228, 'DB040105': 229, 'DB040106': 230, 'DB040107': 231, 'DB040108': 232, 'DB040109': 233, 'DB040110': 234, 'DB040111': 235, 'DB040112': 236, 'DB040113': 237, 'DB040114': 238, 'DB040115': 239, 'DB040116': 240, 'DB040117': 241, 'DB040118': 242, 'DB040119': 243, 'DB040120': 244, 'DB040121': 245, 'DB040122': 246, 'DB040123': 247, 'DB040124': 248, 'DB040125': 249, 'DB040126': 250, 'DB040127': 251, 'DB040128': 252, 'DB040129': 253, 'DB040130': 254, 'DB040131': 255, 'DB040132': 256, 'DB040133': 257, 'DB040134': 258, 'DB040135': 259, 'DB040136': 260, 'DB040137': 261, 'DB040138': 262, 'DB040139': 263, 'DB040140': 264, 'DB040141': 265, 'DB040142': 266, 'DB040143': 267, 'DB040144': 268, 'DB040145': 269, 'DB040146': 270, 'DB040147': 271, 'DB040148': 272, 'DB040149': 273, 'DB040150': 274, 'DB040151': 275, 'DB040152': 276, 'DB040153': 277, 'DB040154': 278, 'DB040155': 279, 'DB040156': 280, 'DB040157': 281, 'DB040158': 282, 'DB040159': 283, 'DB040160': 284, 'DB040161': 285, 'DB040162': 286, 'DB040163': 287, 'DB040164': 288, 'DB040165': 289, 'DB040166': 290, 'DB040167': 291, 'DB040168': 292, 'DB040169': 293, 'DB040170': 294, 'DB040171': 295, 'DB040172': 296, 'DB040173': 297, 'DB040174': 298, 'DB040175': 299, 'DB040176': 300, 'DB040177': 301, 'DB040178': 302, 'DB040179': 303, 'DB040180': 304, 'DB040181': 305, 'DB040182': 306, 'DB040183': 307, 'DB040184': 308, 'DB040185': 309, 'DB040186': 310, 'DB040187': 311, 'DB040188': 312, 'DB040189': 313, 'DB040190': 314, 'DB040191': 315, 'DB040192': 316, 'DB040193': 317, 'DB040194': 318, 'DB040195': 319, 'DB040196': 320, 'DB040197': 321, 'DB040198': 322, 'DB040199': 323, 'DB040200': 324, 'DB040201': 325, 'DB040202': 326, 'DB040203': 327, 'DB040204': 328, 'DB040205': 329, 'DB040206': 330, 'DB040207': 331, 'DB040208': 332, 'DB040209': 333, 'DB040210': 334, 'DB040211': 335, 'DB040212': 336, 'DB040213': 337, 'DB040214': 338, 'DB040215': 339, 'DB040216': 340, 'DB040217': 341, 'DB040218': 342, 'DB040219': 343, 'DB040220': 344, 'DB040221': 345, 'DB040222': 346, 'DB040223': 347, 'DB040224': 348, 'DB040225': 349, 'DB040226': 350, 'DB040227': 351, 'DB040228': 352, 'DB040229': 353, 'DB040230': 354, 'DB040231': 355, 'DB040232': 356, 'DB040233': 357, 'DB040234': 358, 'DB040235': 359, 'DB040236': 360, 'DB040237': 361, 'DB040238': 362, 'DB040239': 363, 'DB040240': 364, 'DB040241': 365, 'DB040242': 366, 'DB040243': 367, 'DB040244': 368, 'DB040245': 369, 'DB040246': 370, 'DB040247': 371, 'DB040248': 372, 'DB040249': 373, 'DB040250': 374, 'DB040251': 375, 'DB040252': 376, 'DB040253': 377, 'DB040254': 378, 'DB040255': 379, 'DB040256': 380, 'DB040257': 381, 'DB040258': 382, 'DB040259': 383, 'DB040260': 384, 'DB040261': 385, 'DB040262': 386, 'DB040263': 387, 'DB040264': 388, 'DB040265': 389, 'DB040266': 390, 'DB040267': 391, 'DB040268': 392, 'DB040269': 393, 'DB040270': 394, 'DB040271': 395, 'DB040272': 396, 'DB040273': 397, 'DB040274': 398, 'DB040275': 399, 'DB040276': 400, 'DB040277': 401, 'DB040278': 402, 'DB040279': 403, 'DB040280': 404, 'DB040281': 405, 'DB040282': 406, 'DB040283': 407, 'DB040284': 408, 'DB040285': 409, 'DB040286': 410, 'DB040287': 411, 'DB040288': 412, 'DB040289': 413, 'DB040290': 414, 'DB040291': 415, 'DB040292': 416, 'DB040293': 417, 'DB040294': 418, 'DB040295': 419, 'DB040296': 420, 'DB040297': 421, 'DB040298': 422, 'DB040299': 423, 'DB040300': 424, 'DB040301': 425, 'DB040302': 426, 'DB040303': 427, 'DB040304': 428, 'DB040305': 429, 'DB040306': 430, 'DB040307': 431, 'DB040308': 432, 'DB040309': 433, 'DB040310': 434, 'DB040311': 435, 'DB040312': 436, 'DB040313': 437, 'DB040314': 438, 'DB040315': 439, 'DB040316': 440, 'DB040317': 441, 'DB040318': 442, 'DB040319': 443, 'DB040320': 444, 'DB040321': 445, 'DB040322': 446, 'DB040323': 447, 'DB040324': 448, 'DB040325': 449, 'DB040326': 450, 'DB040327': 451, 'DB040328': 452, 'DB040329': 453, 'DB040330': 454, 'DB040331': 455, 'DB040332': 456, 'DB040333': 457, 'DB040334': 458, 'DB040335': 459, 'DB040336': 460, 'DB040337': 461, 'DB040338': 462, 'DB040339': 463, 'DB040340': 464, 'DB040341': 465, 'DB040342': 466, 'DB040343': 467, 'DB040344': 468, 'DB040345': 469, 'DB040346': 470, 'DB040347': 471, 'DB040348': 472, 'DB040349': 473, 'DB040350': 474, 'DB040351': 475, 'DB040352': 476, 'DB040353': 477, 'DB040354': 478, 'DB040355': 479, 'DB040356': 480, 'DB040357': 481, 'DB040358': 482, 'DB040359': 483, 'DB040360': 484, 'DB040361': 485, 'DB040362': 486, 'DB040363': 487, 'DB040364': 488, 'DB040365': 489, 'DB040366': 490, 'DB040367': 491, 'DB040368': 492, 'DB040369': 493, 'DB040370': 494, 'DB040371': 495, 'DB040372': 496, 'DB040373': 497, 'DB040374': 498, 'DB040375': 499, 'DB040376': 500, 'DB040377': 501, 'DB040378': 502, 'DB040379': 503, 'DB040380': 504, 'DB040381': 505, 'DB040382': 506, 'DB040383': 507, 'DB040384': 508, 'DB040385': 509, 'DB040386': 510, 'DB040387': 511, 'DB040388': 512, 'DB040389': 513, 'DB040390': 514, 'DB040391': 515, 'DB040392': 516, 'DB040393': 517, 'DB040394': 518, 'DB040395': 519, 'DB040396': 520, 'DB040397': 521, 'DB040398': 522, 'DB040399': 523, 'DB040400': 524, 'DB040401': 525, 'DB040402': 526, 'DB040403': 527, 'DB040404': 528, 'DB040405': 529, 'DB040406': 530, 'DB040407': 531, 'DB040408': 532, 'DB040409': 533, 'DB040410': 534, 'DB040411': 535, 'DB040412': 536, 'DB040413': 537, 'DB040414': 538, 'DB040415': 539, 'DB040416': 540, 'DB040417': 541, 'DB040418': 542, 'DB040419': 543, 'DB040420': 544, 'DB040421': 545, 'DB040422': 546, 'DB040423': 547, 'DB040424': 548, 'DB040425': 549, 'DB040426': 550, 'DB040427': 551, 'DB040428': 552, 'DB040429': 553, 'DB040430': 554, 'DB040431': 555, 'DB040432': 556, 'DB040433': 557, 'DB040434': 558, 'DB040435': 559, 'DB040436': 560, 'DB040437': 561, 'DB040438': 562, 'DB040439': 563, 'DB040440': 564, 'DB040441': 565, 'DB040442': 566, 'DB040443': 567, 'DB040444': 568, 'DB040445': 569, 'DB040446': 570, 'DB040447': 571, 'DB040448': 572, 'DB040449': 573, 'DB040450': 574, 'DB040451': 575, 'DB040452': 576, 'DB040453': 577, 'DB040454': 578, 'DB040455': 579, 'DB040456': 580, 'DB040457': 581, 'DB040458': 582, 'DB040459': 583, 'DB040460': 584, 'DB040461': 585, 'DB040462': 586, 'DB040463': 587, 'DB040464': 588, 'DB040465': 589, 'DB040466': 590, 'DB040467': 591, 'DB040468': 592, 'DB040469': 593, 'DB040470': 594, 'DB040471': 595, 'DB040472': 596, 'DB040473': 597, 'DB040474': 598, 'DB040475': 599, 'DB040476': 600, 'DB040477': 601, 'DB040478': 602, 'DB040479': 603, 'DB040480': 604, 'DB040481': 605, 'DB040482': 606, 'DB040483': 607, 'DB040484': 608, 'DB040485': 609, 'DB040486': 610, 'DB040487': 611, 'DB040488': 612, 'DB040489': 613, 'DB040490': 614, 'DB040491': 615, 'DB040492': 616, 'DB040493': 617, 'DB040494': 618, 'DB040495': 619, 'DB040496': 620, 'DB040497': 621, 'DB040498': 622, 'DB040499': 623, 'DB040500': 624, 'DB040501': 625, 'DB040502': 626, 'DB040503': 627, 'DB040504': 628, 'DB040505': 629, 'DB040506': 630, 'DB040507': 631, 'DB040508': 632, 'DB040509': 633, 'DB040510': 634, 'DB040511': 635, 'DB040512': 636, 'DB040513': 637, 'DB040514': 638, 'DB040515': 639, 'DB040516': 640, 'DB040517': 641, 'DB040518': 642, 'DB040519': 643, 'DB040520': 644, 'DB040521': 645, 'DB040522': 646, 'DB040523': 647, 'DB040524': 648, 'DB040525': 649, 'DB040526': 650, 'DB040527': 651, 'DB040528': 652, 'DB040529': 653, 'DB040530': 654, 'DB040531': 655, 'DB040532': 656, 'DB040533': 657, 'DB040534': 658, 'DB040535': 659, 'DB040536': 660, 'DB040537': 661, 'DB040538': 662, 'DB040539': 663, 'DB040540': 664, 'DB040541': 665, 'DB040542': 666, 'DB040543': 667, 'DB040544': 668, 'DB040545': 669, 'DB040546': 670, 'DB040547': 671, 'DB040548': 672, 'DB040549': 673, 'DB040550': 674, 'DB040551': 675, 'DB040552': 676, 'DB040553': 677, 'DB040554': 678, 'DB040555': 679, 'DB040556': 680, 'DB040557': 681, 'DB040558': 682, 'DB040559': 683, 'DB040560': 684, 'DB040561': 685, 'DB040562': 686, 'DB040563': 687, 'DB040564': 688, 'DB040565': 689, 'DB040566': 690, 'DB040567': 691, 'DB040568': 692, 'DB040569': 693, 'DB040570': 694, 'DB040571': 695, 'DB040572': 696, 'DB040573': 697, 'DB040574': 698, 'DB040575': 699, 'DB040576': 700, 'DB040577': 701, 'DB040578': 702, 'DB040579': 703, 'DB040580': 704, 'DB040581': 705, 'DB040582': 706, 'DB040583': 707, 'DB040584': 708, 'DB040585': 709, 'DB040586': 710, 'DB040587': 711, 'DB040588': 712, 'DB040589': 713, 'DB040590': 714, 'DB040591': 715, 'DB040592': 716, 'DB040593': 717, 'DB040594': 718, 'DB040595': 719, 'DB040596': 720, 'DB040597': 721, 'DB040598': 722, 'DB040599': 723, 'DB040600': 724, 'DB040601': 725, 'DB040602': 726, 'DB040603': 727, 'DB040604': 728, 'DB040605': 729, 'DB040606': 730, 'DB040607': 731, 'DB040608': 732, 'DB040609': 733, 'DB040610': 734, 'DB040611': 735, 'DB040612': 736, 'DB040613': 737, 'DB040614': 738, 'DB040615': 739, 'DB040616': 740, 'DB040617': 741, 'DB040618': 742, 'DB040619': 743, 'DB040620': 744, 'DB040621': 745, 'DB040622': 746, 'DB040623': 747, 'DB040624': 748, 'DB040625': 749, 'DB040626': 750, 'DB040627': 751, 'DB040628': 752, 'DB040629': 753, 'DB040630': 754, 'DB040631': 755, 'DB040632': 756, 'DB040633': 757, 'DB040634': 758, 'DB040635': 759, 'DB040636': 760, 'DB040637': 761, 'DB040638': 762, 'DB040639': 763, 'DB040640': 764, 'DB040641': 765, 'DB040642': 766, 'DB040643': 767, 'DB040644': 768, 'DB040645': 769, 'DB040646': 770, 'DB040647': 771, 'DB040648': 772, 'DB040649': 773, 'DB040650': 774, 'DB040651': 775, 'DB040652': 776, 'DB040653': 777, 'DB040654': 778, 'DB040655': 779, 'DB040656': 780, 'DB040657': 781, 'DB040658': 782, 'DB040659': 783, 'DB040660': 784, 'DB040661': 785, 'DB040662': 786, 'DB040663': 787, 'DB040664': 788, 'DB040665': 789, 'DB040666': 790, 'DB040667': 791, 'DB040668': 792, 'DB040669': 793, 'DB040670': 794, 'DB040671': 795, 'DB040672': 796, 'DB040673': 797, 'DB040674': 798, 'DB040675': 799, 'DB040676': 800, 'DB040677': 801, 'DB040678': 802, 'DB040679': 803, 'DB040680': 804, 'DB040681': 805, 'DB040682': 806, 'DB040683': 807, 'DB040684': 808, 'DB040685': 809, 'DB040686': 810, 'DB040687': 811, 'DB040688': 812, 'DB040689': 813, 'DB040690': 814, 'DB040691': 815, 'DB040692': 816, 'DB040693': 817, 'DB040694': 818, 'DB040695': 819, 'DB040696': 820, 'DB040697': 821, 'DB040698': 822, 'DB040699': 823, 'DB040700': 824, 'DB040701': 825, 'DB040702': 826, 'DB040703': 827, 'DB040704': 828, 'DB040705': 829, 'DB040706': 830, 'DB040707': 831, 'DB040708': 832, 'DB040709': 833, 'DB040710': 834, 'DB040711': 835, 'DB040712': 836, 'DB040713': 837, 'DB040714': 838, 'DB040715': 839, 'DB040716': 840, 'DB040717': 841, 'DB040718': 842, 'DB040719': 843, 'DB040720': 844, 'DB040721': 845, 'DB040722': 846, 'DB040723': 847, 'DB040724': 848, 'DB040725': 849, 'DB040726': 850, 'DB040727': 851, 'DB040728': 852, 'DB040729': 853, 'DB040730': 854, 'DB040731': 855, 'DB040732': 856, 'DB040733': 857, 'DB040734': 858, 'DB040735': 859, 'DB040736': 860, 'DB040737': 861, 'DB040738': 862, 'DB040739': 863, 'DB040740': 864, 'DB040741': 865, 'DB040742': 866, 'DB040743': 867, 'DB040744': 868, 'DB040745': 869, 'DB040746': 870, 'DB040747': 871, 'DB040748': 872, 'DB040749': 873, 'DB040750': 874, 'DB040751': 875, 'DB040752': 876, 'DB040753': 877, 'DB040754': 878, 'DB040755': 879, 'DB040756': 880, 'DB040757': 881, 'DB040758': 882, 'DB040759': 883, 'DB040760': 884, 'DB040761': 885, 'DB040762': 886, 'DB040763': 887, 'DB040764': 888, 'DB040765': 889, 'DB040766': 890, 'DB040767': 891, 'DB040768': 892, 'DB040769': 893, 'DB040770': 894, 'DB040771': 895, 'DB040772': 896, 'DB040773': 897, 'DB040774': 898, 'DB040775': 899, 'DB040776': 900, 'DB040777': 901, 'DB040778': 902, 'DB040779': 903, 'DB040780':
```

Then by using the `map()` function, a new column of “DrugBankIDLabel” can be generated and its values are integers.

```
data["DrugBankIDLabel"] = data["DrugBankID"].map(dict1)
data.head()
```

DrugBankID	Name	Type	UniProtID	UniProtName	DrugBankIDLabel	
0	DB00001	Lepirudin	BiotechDrug	P00734	Prothrombin	593
1	DB00002	Cetuximab	BiotechDrug	P00533	Epidermal growth factor receptor	6639
2	DB00002	Cetuximab	BiotechDrug	O75015	Low affinity immunoglobulin gamma Fc region re...	6639
3	DB00002	Cetuximab	BiotechDrug	P00736	Complement C1r subcomponent	6639
4	DB00002	Cetuximab	BiotechDrug	P02745	Complement C1q subcomponent subunit A	6639

Similarly, a new column of “UniProtIDLabel” is generated and its values are integers.

```
data["UniProtIDLabel"] = data["UniProtID"].map(dict2)
data.head()
```

DrugBankID	Name	Type	UniProtID	UniProtName	DrugBankIDLabel	UniProtIDLabel	
0	DB00001	Lepirudin	BiotechDrug	P00734	Prothrombin	593	1299
1	DB00002	Cetuximab	BiotechDrug	P00533	Epidermal growth factor receptor	6639	3561
2	DB00002	Cetuximab	BiotechDrug	O75015	Low affinity immunoglobulin gamma Fc region re...	6639	4525
3	DB00002	Cetuximab	BiotechDrug	P00736	Complement C1r subcomponent	6639	3283
4	DB00002	Cetuximab	BiotechDrug	P02745	Complement C1q subcomponent subunit A	6639	2352

Because all the drug-target pairs from DrugBank are proved to be interacting with each other, the interaction probability is 100%. We generated another new column called “Interaction” and put 1 as the values for all these reported drug-target pairs. We call these positive samples. The data set with only positive samples looks like below (20061 rows × 3 columns):

	DrugBankIDLabel	UniProtIDLabel	Interaction
0	593	1299	1
1	6639	3561	1
2	6639	4525	1
3	6639	3283	1
4	6639	2352	1
5	6639	2487	1
6	6639	2212	1
7	6639	2450	1
8	6639	2193	1
9	6639	432	1
10	6639	114	1
11	6639	1277	1
12	6639	266	1
13	1146	2731	1
14	1146	3575	1
15	1146	2003	1
16	2515	971	1
17	2515	931	1
18	2515	432	1
19	2515	2450	1

In order to train a good model, we have to generate negative samples. In our case, the negative sample means a drug-target pair with an interaction probability of 0. The Drugbank database only provides positive samples while there are no experimentally-validated negative samples. So we

can only assume some drug-target pairs have 0 probability to be interacting with each other although they must not be real negative samples. The way we generate the negative samples is to randomly generate a drug-target pair and then verify if it exists in the positive sample set, if yes, then re-generate a new pair, if no, this pair is saved, and then go back the for loop to generate another pair until the loop ends. The below code iterates 30000 times and “data1” finally includes 20061 positive samples and 29992 negative samples.

```
import random
for i in range(30000):
    a=random.randint(1,7370)
    b=random.randint(1,4763)
    if (a,b) in dict3:
        continue
    data1.loc[20061+i]=[a,b,0]
data1
```

	DrugBankIDLabel	UniProtIDLabel	Interaction
0	593	1299	1
1	6639	3561	1
2	6639	4525	1
3	6639	3283	1
4	6639	2352	1
5	6639	2487	1
6	6639	2212	1
7	6639	2450	1
8	6639	2193	1

The above data set was split into train set and test set by:

```
(train, test) = train_test_split(data, train_size = 0.8)
```

We have to assign a value for the below parameters of the PMF model:

Parameter	Meaning
n_user	It is the total number of drugs in our case. The value is 7370.
n_item	It is the total number of targets in our case. The value is 4763.
n_feature	It tells how many features can uniquely describe a drug and a target. It composes an embedding vector for each drug and each target. Intrinsically, the inner product of the drug's embedding vector and the target's embedding vector gives the interaction probability of this drug and this target. We plan to run the PMF model a few times to find an optimal n_feature value.
n_iters	The number of iterations. The higher value of n_iters is supposed to give a better prediction but the running time will be increasing drastically.
max_rating	In our case, it is 1.
min_rating	In our case, it is 0.

We firstly tried this combination of parameters: n_user=7370, n_item=4763, n_feature=30, n_iters=10, max_rating=1, min_rating=0. And we obtained below results:

```

print("training pmf...")
pmf = PMF(n_user=n_user, n_item=n_items, n_feature = n_feat, epsilon = 5, max_rating=1, min_rating=0)
pmf.fit(train, n_iters=eval_iters)

training pmf...
<recommend.pmf.PMF at 0x267290190f0>

```

```
from recommend.utils.evaluation import RMSE

train_preds = pmf.predict(train[:, :2])
train_rmse = RMSE(train_preds, train[:, 2])

test_preds = pmf.predict(test[:, :2])
test_rmse = RMSE(test_preds, test[:, 2])

print('Train RMSE: ', train_rmse)
print('Test RMSE: ', test_rmse)
print()
print('Train Preds: ')
print(train_preds)
print()
print("Test Preds: ")
print(test_preds)

print("Test min: ", test_preds.min())
print('Test max: ', test_preds.max())
print("Train min: ", train_preds.min())
print('Train max: ', train_preds.max())
```

Train RMSE: 0.49113256485741547

Test RMSE: 0.491773108704835

Train Preds:

[0.48347627 0.46815479 0.54032976 ... 0.48009259 0.47584957 0.48551968]

Test Preds:

[0.49458976 0.48517228 0.46688036 ... 0.46846345 0.504811 0.47957349]

Test min: 0.4341419607311734

Test max: 0.5835814951246646

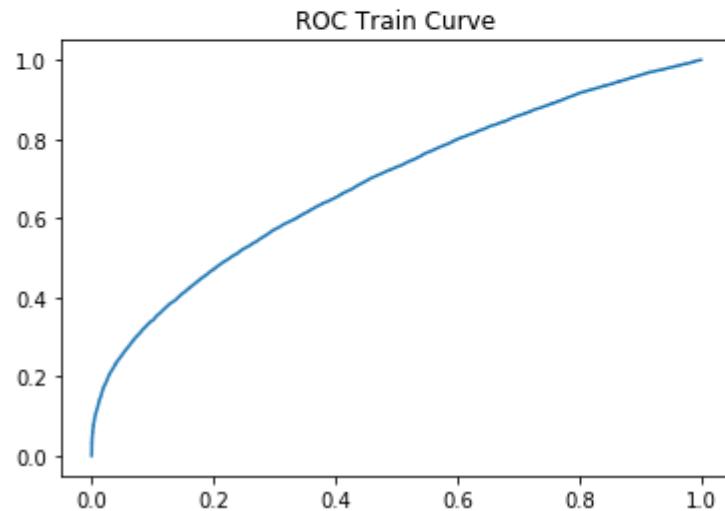
Train min: 0.42807871183165375

Train max: 0.591980499102021

```
(fpr_train, tpr_train, thresholds_train) = roc_curve(train[:,2],train_preds)
auc_train = roc_auc_score(train[:,2],train_preds)
```

```
plt.plot(fpr_train,tpr_train)
plt.title('ROC Train Curve')
print("AUC Train: ", auc_train)
```

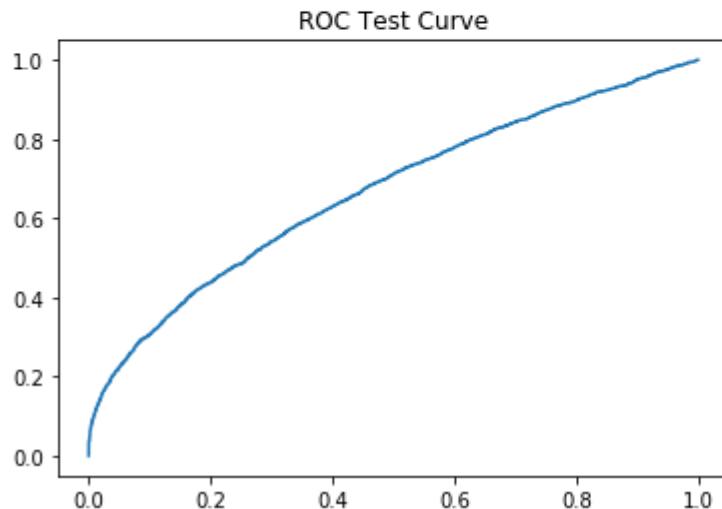
AUC Train: 0.6877667795585204



```
(fpr, tpr, thresholds) = roc_curve(test[:,2],test_preds)
auc_test = roc_auc_score(test[:,2],test_preds)
```

```
plt.plot(fpr,tpr)
plt.title("ROC Test Curve")
print('AUC Test: ', auc_test)
```

AUC Test: 0.6664035608358009



The results are summarized in below table:

Performance Indicator	Value
RMSE for train set	0.4911
RMSE for test set	0.4918
AUC score for train set	0.6878
AUC score for test set	0.6664

Results for: n_user=7370, n_item=4763, n_feature=30, n_iters=10, max_rating=1, min_rating=0

Similarly, we ran the model using different combinations of parameters and the results are listed below:

Performance Indicator	Value
RMSE for train set	0.5119
RMSE for test set	0.5155
AUC score for train set	0.7331
AUC score for test set	0.6875

Results for: n_user=7370, n_item=4763, n_feature=70, n_iters=10, max_rating=1, min_rating=0

Performance Indicator	Value
RMSE for train set	0.4643
RMSE for test set	0.4703
AUC score for train set	0.8773
AUC score for test set	0.7966

Results for: n_user=7370, n_item=4763, n_feature=30, n_iters=50, max_rating=1, min_rating=0

Performance Indicator	Value
RMSE for train set	0.4443
RMSE for test set	0.4585
AUC score for train set	0.9078
AUC score for test set	0.8010

Results for: n_user=7370, n_item=4763, n_feature=30, n_iters=100, max_rating=1, min_rating=0

Comparing the above four results, the current optimal result occurred at n_feature=30, n_iters=100. The predicted data was exported to “PMF Generated data_30000_fea30_ite100.csv”. The current model can be used to predict the interaction probability for any drug-target pair but with limited accuracy. Taking DrugBankIDLabel #1 for example, we paired it with all the 4763 targets and predicted the interaction probability of all these pairs. Below is the screenshot of the result:

DrugBankIDLabel	UniProtIDLabel	Predicted interaction
1	1	0.432771164
1	2	0.447979408
1	3	0.428092253
1	4	0.518644897
1	5	0.429846909
1	6	0.425531057
1	7	0.429120209
1	8	0.463735628
1	9	0.454380813
1	10	0.430469087
1	11	0.416836863
1	12	0.45221196
1	13	0.424542773
1	14	0.562239136
1	15	0.457691837
1	16	0.427313679
1	17	0.439664403
1	18	0.44841151
1	19	0.457195833
1	20	0.444873125

Q3:

It has been reported that the quality of the negative data set affects the prediction results [3]. The chemical and structural similarity has been incorporated into the procedure of generating

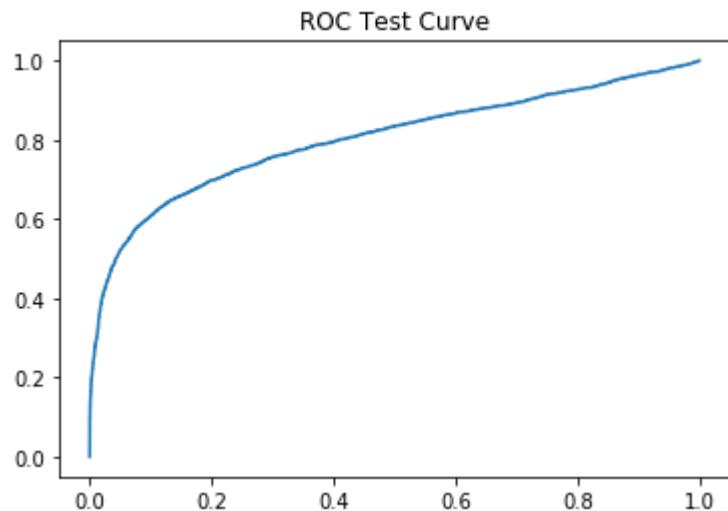
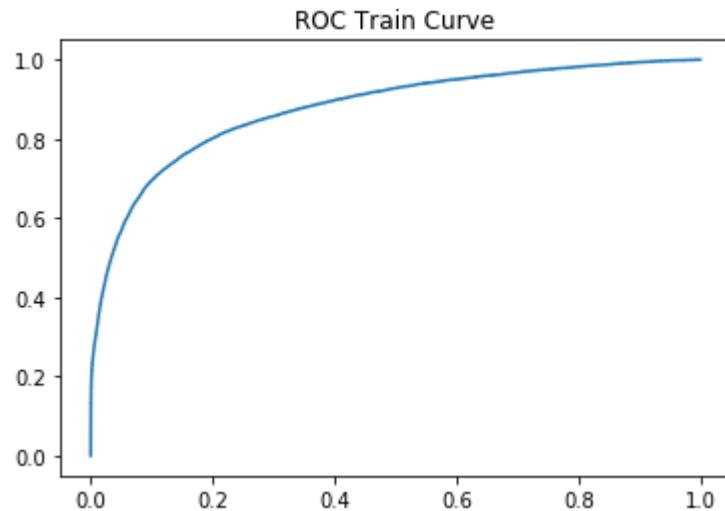
the negative data set [2]. In our study, we did randomly pairing of the drug and the target and assume their interaction probability to be 0. And we varied the number of negative samples to form different data sets aiming to see how the number of negative samples affects the prediction performance.

We generated another data set with 20061 positive samples and 79955 negative samples. We split the data set by using “(train, test) = train_test_split(data, train_size = 0.8) and used n_feature=30, n_iters=100 to fit the PMF model on the train data. The results are summarized as below:

Performance Indicator	Value
RMSE for train set	0.3692
RMSE for test set	0.3771
AUC score for train set	0.8760
AUC score for test set	0.8054

Results for: n_user=7370, n_item=4763, n_feature=30, n_iters=100, max_rating=1, min_rating=0

The ROC curves as below:



We can see that when the number of the negative samples increases from 29992 to 79955, the RMSE is lowered by about 17.8% and the AUC score almost stays the same. The reduction in RMSE might be because more negative samples bring more zero(s) into the calculation. So it looks to us that there is no improvement in the prediction accuracy when there are more negative samples.

Problem Report

The predicted interaction probability for the positive samples is supposed to be very close to 1 but it is actually much smaller than 1 as shown below. And as we have seen in the tables above, the RMSE and AUC score still have some room to improve. The future work will have to take the chemical similarity into consideration when generating the negative samples. For example, given a known interaction of drug i and target j , for a drug k , the more difference between drug i and drug k , the less possibility that drug k interacts target j . We expect to see improved prediction accuracy after the chemical similarity is taken into account.

	DrugBankIDLabel	UniProtIDLabel	Interaction	Predicted interaction
0	592	1298	1	0.701755798
1	6638	3560	1	0.700387909
2	6638	4524	1	0.591572019
3	6638	3282	1	0.659930284
4	6638	2351	1	0.598876777
5	6638	2486	1	0.630448546
6	6638	2211	1	0.61587401
7	6638	2449	1	0.620322638
8	6638	2192	1	0.580231696
9	6638	431	1	0.610358415
10	6638	113	1	0.570520987
11	6638	1276	1	0.636257339
12	6638	265	1	0.594370809
13	1145	2730	1	0.486199189
14	1145	3574	1	0.440785964
15	1145	2002	1	0.461842285
16	2514	970	1	0.656506231
17	2514	930	1	0.476426109
18	2514	431	1	0.592702431
19	2514	2449	1	0.59423371

50032	4509	2207	0	0.45115946
50033	5613	3410	0	0.448661049
50034	5817	1219	0	0.453267948
50035	1707	1696	0	0.501528024
50036	3961	1882	0	0.452753702
50037	5268	706	0	0.429287766
50038	4674	1263	0	0.437456006
50039	384	2921	0	0.455278827
50040	6756	1040	0	0.474904974
50041	6599	2538	0	0.450695837
50042	5438	625	0	0.444811952
50043	3922	2541	0	0.437489241
50044	3655	1630	0	0.429225442
50045	3816	4045	0	0.462735039
50046	57	4336	0	0.446154686
50047	49	2174	0	0.450375719
50048	8	990	0	0.427661522
50049	3198	2865	0	0.436243197
50050	142	781	0	0.435927561
50051	4059	3580	0	0.474465739
50052	6014	3479	0	0.482048911

Screenshots of the partial prediction results for the data set with 29992 negative samples and n_feature=30, n_iters=100.

References

1. Ruslan Salakhutdinov and Andriy Mnih. Probabilistic matrix factorization. In Advances in Neural Information Processing Systems; NIPS, 2007 20: 1257-1264
2. Implementation of PMF:
<https://github.com/chyikwei/recommend/tree/master/recommend>
3. Zhanzhan Cheng, Kai Huang, Yang Wang, Hui Liu, Jihong Guan and Shuigeng Zhou. Selecting high-quality negative samples for effectively predicting protein-RNA interactions. BMC Systems Biology. 2017 11 (Suppl 2) :9