

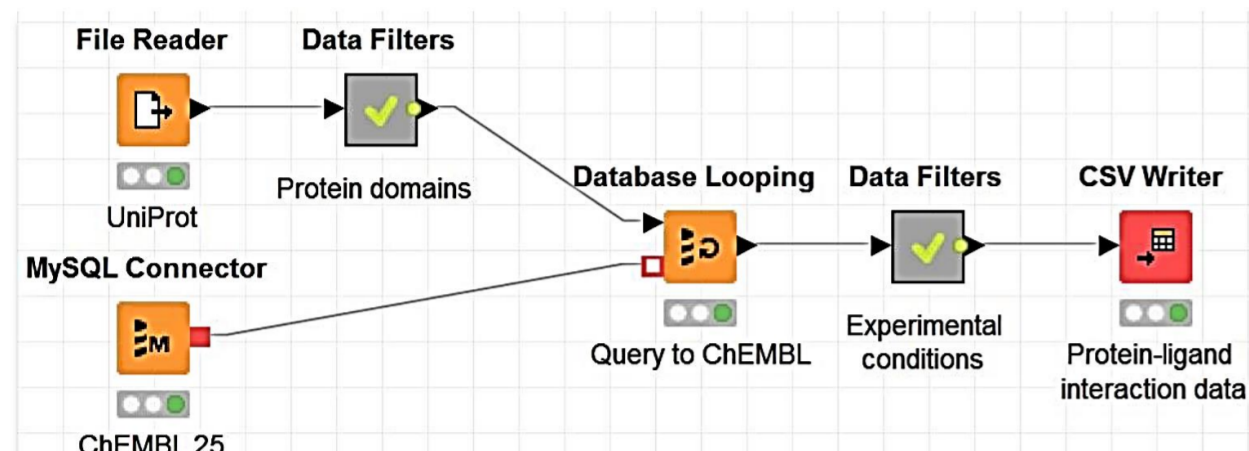
Введение

Компьютерное прогнозирование взаимодействий «белок-лиганд» - важный этап поиска фармакологических веществ. Полученные знания могут также служить для установления оптимальных направлений лекарственной терапии. В связи со значительной неполнотой экспериментальных данных по биологической активности фармакологических веществ при поиске белков-мишеней применяют методы машинного обучения. Цель работы заключается в разработке метода, который позволит эффективно прогнозировать взаимодействия для разнообразных групп белков и химических соединений.

Методы

1. Обучающие данные

Для извлечения информации о белках-мишенях и взаимодействующих с ними лигандами нами был разработан автоматизированный конвейер на платформе KNIME (https://www.knime.com).



Источником информации служила база данных ChEMBL. Отобраны результаты исследований in vitro, с указанием показателей аффинности IC50, Kd или Ki. Проведена фильтрация мутантных форм белков (таб. 1).

Таблица 1. Характеристики собранных данных

Группа белков	Activity cutoff (µM)	Ligand	Target
GPCR	IC50	1	15025
		10	17488
	Kd	1	557
		10	612
		1	32490
		10	35888
Ion Channel (LG*)	IC50	1	1840
		10	1980
	Kd	1	13
		10	15
		1	1018
		10	1165
Ion Channel (VG**)	IC50	1	4809
		10	8472
	Kd	1	7
		10	8
		1	976
		10	1485
Nuclear receptors	IC50	1	3832
		10	4827
	Kd	1	143
		10	168
		1	1134
		10	1238
Protein kinases	IC50	1	33078
		10	38795
	Kd	1	489
		10	574
		1	4077
		10	4410

2. Алгоритм

Использовался метод PASS, который был разработан для прогноза биологической активности лекарственно подобных органических соединений [1]. В данном методе используются оригинальные дескрипторы многоуровневых окрестностей атомов (MNA) (рис. 1), которые отражают локальные особенности взаимодействия лиганд-мишень. Метод основан на наивном байесовском классификаторе. Для каждой пары мишень-лиганд программа PASS рассчитывает оценки Pa и Pi, которые представляют собой вероятности взаимодействия и отсутствия взаимодействия между целевым белком и лигандом, соответственно.

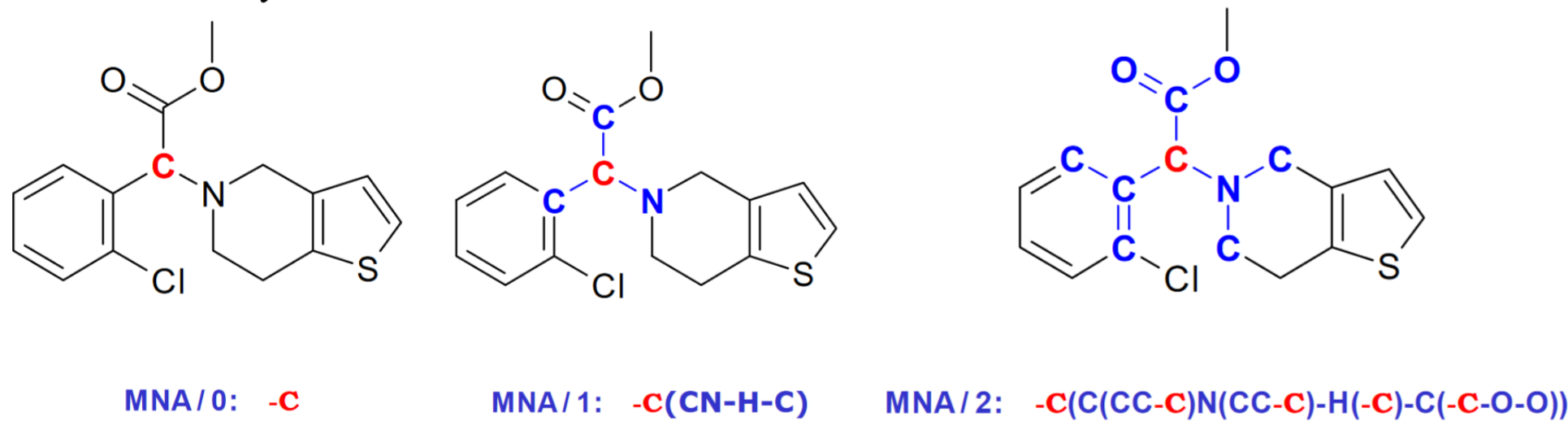


Рисунок 1. Пример генерации MNA-дескрипторов нулевого (MNA/0), первого (MNA/1) и второго (MNA/2) уровней.

Для оценки последовательностей использовалась программа SPrOS [2], в которой реализован алгоритм, основанный на оценках локального сходства аминокислотных последовательностей. Аминокислотная последовательность тестовой белка (Q) сопоставляется с первичными структурами каждого из белков обучающей выборки (K) путем серии сдвигов.

```

AANRDPSPQFPDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 2
ANRDPSPQFPDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 1
NRDPSPQFPDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 1
RDPSPQFPDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 0
DPSQFPDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 1
PSQFPDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 2
SQFPDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 1
QFPDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 1
FPDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 2
DPDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 0
PDPHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 1
PHRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 0
HRRFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 9
RFDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 0
FDVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 3
DVTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 1
VTRDTRGHLSPFGQGIHFCMGRPLAKLEGEVALR 1
    
```

GTAINKPLSEKMLFGMGKRRICIGEVLAKEWEIFLFLAILLQQLEFSV 9
Последовательность Q

$$R_{ih} = \sum_{j=i}^{j+F-1} \text{sim}(q_j, k_{j+h})$$

K – последовательность обучающей выборки
Q – тестовая последовательность
i – позиция в тестовой последовательности Q
F – длина сегментов сопоставления
sim(q, k) – подобие аминокислотных остатков в совмещенных позициях (идентичность)
h – величина сдвига между последовательностями Q и K

$$S_p = \max_{h,i} R_{ih}, p - F < i \leq p$$

Каждой позиции p присваивается наибольшее значение из оценок Ri

Оценка принадлежности белка к классам лигандной специфичности

$$t_p = \frac{\sum_{k=1}^N S_{pk} \times [a_k(C) - b_k(C)]}{\sum_{k=1}^N S_{pk} \times [a_k(C) + b_k(C)]}$$

C – лиганд обучающей выборки
tp – интегральная оценка для каждой позиции тестовой последовательности
Spk – оценка позиции p при сопоставлении с k-той последовательностью обучающей выборки
ak(C), bk(C) – коэффициенты принадлежности белка k к классу лигандной специфичности C и к его дополнению, соответственно

$$t = \sin\left[\frac{1}{m} \sum_{p=1}^m \arcsin(t_p)\right]$$

Усреднение оценок tp для всех позиций m тестовой последовательности Q

$$t_0 = \frac{\sum_{k=1}^N [a_k(C) - b_k(C)]}{\sum_{k=1}^N [a_k(C) + b_k(C)]}$$

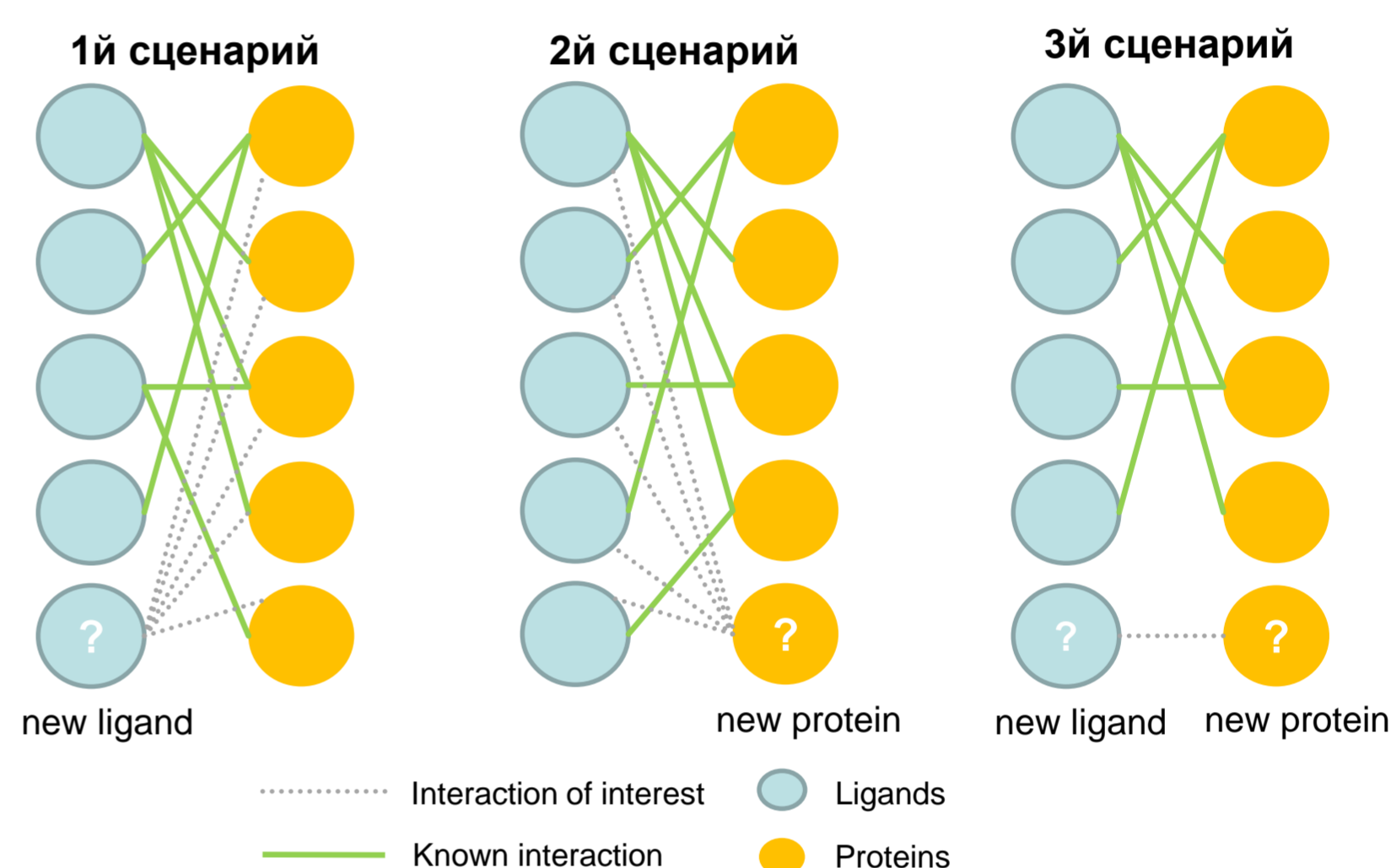
Для учета представительности класса в обучающей выборке рассчитывается t0

$$B(C) = \frac{t - t_0}{1 - t_0}$$

Принадлежность белка Q к классу лигандной специфичности C оценивается с помощью B-статистики

Результаты

Тестирование программы осуществлялось на четырех наборах данных: рецепторы, связанные с G-белком (GPCR), протеинкиназы, ядерные рецепторы, ионные каналы (таб. 1). Тестирование осуществлялось в трех различных сценариях исследований (рис. 2).



При первом сценарии осуществляется прогноз связывания с белками, для которых уже известны лиганды. При втором – для нового белкам-мишени предсказывались лиганды, для которых уже известны какие-либо белки-мишени. Третий сценарий – для нового белка-мишени необходимо предсказать с новым лигандом.

Таблица 2. Результаты тестирования.

Группа белков	Показатель взаимодействия	Порог активности (µmol)	Первый сценарий (AUC)	Второй сценарий (AUC)	Третий сценарий (AUC)
GPCR	IC50	1	0,9864	0,968	0,904
		10	0,981	0,953	0,910
	Kd	1	0,9788	0,875	0,805
		10	0,9809	0,962	0,882
		1	0,9868	0,976	0,901
		10	0,9835	0,977	0,918
Протеинкиназы	IC50	1	0,9625	0,924	0,857
		10	0,954	0,902	0,838
	Kd	1	0,8079	0,790	0,655
		10	0,8026	0,797	0,650
		1	0,9797	0,956	0,893
		10	0,9806	0,937	0,856
Ion channel (ligand-gated)	IC50	1	0,9875	1,000	1,000
		10	0,9857	0,979	0,957
	Kd	1	0,9912	0,973	0,963
		10	0,9857	0,985	0,932
		10	0,9677	0,898	0,839
		1	0,9843	0,973	0,852
Ion channel (voltage-gated)	IC50	1	0,9811	0,988	0,943
		10	0,9811	0,988	0,943
	Kd	1	0,9686	1,000	0,961
		10	0,9715	1,000	0,972
		1	0,9926	0,988	0,976
		10	0,9955	0,987	0,989

Заключение

Таким образом, разработанный подход позволяет эффективно прогнозировать связывание белков с низкомолекулярными органическими соединениями. Высокая точность прогноза получена во всех сценариях прогноза, даже для пар белок-лиганд, когда ни для лиганда, ни для белка мишени ничего не известно. Наш метод применим для групп белков с различной степенью корреляции между филогенетическими отношениями и специфичностью к низкомолекулярным лигандам.

Работа выполнена при поддержке гранта РФФИ № 19-015-00374.

Ссылки:

- Filimonov DA, Lagunin AA, Glorizova TA, Rudik AV, Druzhilovskii DS, Pogodin PV, Poroikov VV, Prediction of the biological activity spectra of organic compounds using the PASS online web resource, *Chemistry of Heterocycl Compd* **50**: 444-457, 2014.
- Karasev DA, Sobolev BN, Lagunin AA, Filimonov DA, Poroikov VV, Prediction of Protein–ligand Interaction Based on Sequence Similarity and Ligand Structural Features, *Int. J. Mol. Sci.* **2020**, *21*(21), 8152