

Predicting the activity of protein-ligand complexes

Lukas Fallmann



BACHELORARBEIT

eingereicht am
Fachhochschul-Bachelorstudiengang

Medizin- und Bioinformatik

in Hagenberg

im Juni 2023

Advisor:

Micha Johannes Birklbauer, M.Sc.

© Copyright 2023 Lukas Fallmann

This work is published under the conditions of the Creative Commons License *Attribution-NonCommercial-NoDerivatives 4.0 International* (CC BY-NC-ND 4.0)—see <https://creativecommons.org/licenses/by-nc-nd/4.0/>.

Declaration

I hereby declare and confirm that this thesis is entirely the result of my own original work. Where other sources of information have been used, they have been indicated as such and properly acknowledged. I further declare that this or similar work has not been submitted for credit elsewhere. This printed copy is identical to the submitted electronic version.

Hagenberg, June 27, 2023

Lukas Fallmann

Contents

Declaration	iv
Preface	vii
Abstract	viii
Kurzfassung	ix
1 Introduction	1
1.1 Overview of the topic	2
1.2 Why is it important	2
1.3 What are current methods that are prominently used	2
1.4 Machine Learning in drug design / activity prediction	2
1.5 Interactions	2
1.6 Goals	2
2 Methods	3
2.1 Data description	3
2.2 Data partitioning	3
2.3 Machine Learning approaches	3
2.4 quality metrics	3
2.5 hyperparameter search	3
2.6 feature engineering	3
3 Results	4
3.1 Results of the models, mostly tables but for best performing models in each category also do some plots like confusion matrices, AUC plots if available, bar plots for comparison between models	4
4 Discussion	5
4.1 Performance	5
4.2 Improvements	5
A Technical Details	6
B Supplementary Materials	7
B.1 PDF Files	7

Contents	vi
B.2 Media Files	7
B.3 Online Sources (PDF Captures)	7
C Questionnaire	8
D LaTeX Source Code	9
References	10
Literature	10
Online sources	10

Preface

Abstract

This should be a 1-page (maximum) summary of your work in English.

Kurzfassung

An dieser Stelle steht eine Zusammenfassung der Arbeit, Umfang max. 1 Seite. ...

Chapter 1

Introduction

The discovery of new drugs or any chemically active compounds for that matter is an expensive and time-consuming process. It has been estimated, that it takes about 14 Years from the initial discovery of a promising new compound to the release of a marketable drug[7]. In addition to that the price of this drug-discovery circle ranges up to 800 Million Dollars[2]. All techniques which aim to improve the efficiency of drug discovery can be generalized as one of two methods. These two are called High-throughput-screening (HTS) and virtual screening (VS).

When using an HTS-approach there are many compounds which are tested against some type of target protein. During testing, it is measured whether a certain compound biochemically interacts with a protein. Those interacting combinations are considered active and are marked by researchers as hits. To improve the performance of HTS there are a number of factors to consider. Through miniaturization, it is possible to investigate more compounds at the same time. With a higher throughput quality-control is more time-consuming and leads to an overall more expensive process. For this reason HTS is most efficient, when analyzing a small set of compounds as the technology is not suitable for large datasets[6].

In contrast to the in vitro approach of HTS, VS is a theoretical in silico approach. To save resources in the laboratory the activity of certain compounds is predicted using a preexisting library of small molecules. The activity can be predicted using the ligands of a compound and their respective binding sites or the 3D structure of a compound. The Key idea behind the ligand based approach (LBVS) is that similar compounds have similar chemical properties. Therefore, the goal of LBVS is to find molecules which have similar or identical chemical properties as the sample compound[3]. Structure-based VS uses the 3D structure of a compound to predict which molecules from the dataset will bind to the provided sample. Each molecule of a certain database subset is fitted (docked) to the sample. Hereby it is important to differentiate between rigid and flexible docking.

In rigid docking the dataset sample is rotated and translated in a six-dimensional space in order to fit the sample protein. For each fitted molecule a score is calculated based on how well the molecule fits to the sample[4]. Although this algorithm often predicts actual possible binding sites and bound proteins, there is no guarantee that this compound will actually bind in vitro. Therefore, predicted interactions should be seen as a hypothesis. Still rigid docking provides a great baseline at a comparatively low cost[3]. The low

accuracy of rigid docking is due to the nature of biochemical substances as samples in a database can only provide a snapshot of a sample. With flexible docking it is possible to simulate moving binding sites. The flexibility can be introduced at different stages. Implicit flexibility is achieved by smoothing protein surfaces and therefore allowing room for interpretation when docking. Cross- or Ensemble docking can be done by repeating the docking process with different conformations. Explicit flexibility is reached through allowing side-chain flexibility. Most commonly utilized is the approach where the ligand is flexible, and the receptor is rigid. Even though this approach does provide better more accurate results it takes considerably longer to compute[8].

Regardless of the docking type the score decides which pose between a protein and a ligand is most likely to exist. In addition to that, the score also determines whether a protein-ligand complex is considered active. There are a lot of different scoring functions which can be grouped into four categories: physics-based, empirical, knowledge-based, and machine learning-based[5].

The focus of this work is on implementing a machine-learning based scoring approach. Machine-learning based scoring functions work by training on pre-classified data and finding the best model for predicting future data. To accurately and efficiently train a model crucial binding sites need to be identified beforehand. The basis of this thesis is the master thesis of Birklbauer Micha[1]. In his thesis a selection of eleven proteins from the directory of useful decoys have been selected to be analyzed. For the selected proteins all possible interactions have been analyzed. Based on the interaction-data a few basic scoring functions have been implemented. The direct result of this thesis are proteins and the frequency of their interactions.

Since this work aims to implement different machine learning algorithms for use in scoring functions the state of the art is described in the following.

1.1 Overview of the topic

1.2 Why is it important

1.3 What are current methods that are prominently used

1.4 Machine Learning in drug design / activity prediction

1.5 Interactions

1.6 Goals

Chapter 2

Methods

- 2.1 Data description
- 2.2 Data partitioning
- 2.3 Machine Learning approaches
- 2.4 quality metrics
- 2.5 hyperparameter search
- 2.6 feature engineering

Chapter 3

Results

- 3.1 Results of the models, mostly tables but for best performing models in each category also do some plots like confusion matrices, AUC plots if available, bar plots for comparison between models

Chapter 4

Discussion

4.1 Performance

4.2 Improvements

Appendix A

Technical Details

Appendix B

Supplementary Materials

List of supplementary data submitted to the degree-granting institution for archival storage (in ZIP format).

B.1 PDF Files

Path: /

thesis.pdf Master/Bachelor thesis (complete document)

B.2 Media Files

Path: /media

*.ai, *.pdf Adobe Illustrator files

*.jpg, *.png raster images

*.mp3 audio files

*.mp4 video files

B.3 Online Sources (PDF Captures)

Path: /online-sources

Reliquienschrein-Wikipedia.pdf [9]

Appendix C

Questionnaire

Appendix D

LaTeX Source Code

References

Literature

- [1] Micha Johannes Birklbauer. “Automatic identification of important interactions and interaction-frequency-based scoring in protein-ligand complexes”. MA thesis. FH Hagenberg, Aug. 31, 2021 (cit. on p. 2).
- [2] Joseph A. DiMasi, Ronald W. Hansen, and Henry G. Grabowski. “The price of innovation: new estimates of drug development costs”. eng. *Journal of Health Economics* 22.2 (Mar. 2003), pp. 151–185. DOI: 10.1016/S0167-6296(02)00126-1 (cit. on p. 1).
- [3] Aleix Gimeno et al. “The Light and Dark Sides of Virtual Screening: What Is There to Know?”. *International Journal of Molecular Sciences* 20.6 (Mar. 2019), p. 1375. DOI: 10.3390/ijms20061375. (Visited on 02/29/2024) (cit. on p. 1).
- [4] A. Lavecchia and C. Di Giovanni. “Virtual screening strategies in drug discovery: a critical review”. eng. *Current Medicinal Chemistry* 20.23 (2013), pp. 2839–2860. DOI: 10.2174/09298673113209990001 (cit. on p. 1).
- [5] Jin Li, Ailing Fu, and Le Zhang. “An Overview of Scoring Functions Used for Protein–Ligand Interactions in Molecular Docking”. en. *Interdisciplinary Sciences: Computational Life Sciences* 11.2 (June 2019), pp. 320–328. DOI: 10.1007/s12539-019-00327-w. (Visited on 02/29/2024) (cit. on p. 2).
- [6] Lorenz M. Mayr and Peter Fuerst. “The Future of High-Throughput Screening”. *SLAS Discovery* 13.6 (July 2008), pp. 443–448. DOI: 10.1177/1087057108319644. (Visited on 02/29/2024) (cit. on p. 1).
- [7] S. Myers and A. Baker. “Drug discovery—an operating model for a new era”. eng. *Nature Biotechnology* 19.8 (Aug. 2001), pp. 727–730. DOI: 10.1038/90765 (cit. on p. 1).
- [8] Nataraj S. Pagadala, Khajamohiddin Syed, and Jack Tuszynski. “Software for molecular docking: a review”. *Biophysical Reviews* 9.2 (Jan. 2017), pp. 91–102. DOI: 10.1007/s12551-016-0247-1. (Visited on 02/29/2024) (cit. on p. 2).

Online sources

- [9] *Reliquienschrein*. Aug. 29, 2022. URL: <https://de.wikipedia.org/wiki/Reliquienschrein> (visited on 02/11/2023).

Check Final Print Size

— Check final print size! —



— Remove this page after printing! —