

# 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

<b>Declaration</b>	<b>iv</b>
<b>Preface</b>	<b>vii</b>
<b>Abstract</b>	<b>viii</b>
<b>Kurzfassung</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Machine Learning in drug design and activity prediction !WIP! . . . . .	3
1.2 Goals . . . . .	3
<b>2 Methods</b>	<b>4</b>
2.1 Data description . . . . .	4
2.1.1 Proteins . . . . .	4
2.1.2 Interactions → !WIP! ask Micha regarding grade of detail . . . .	5
2.1.3 Data origin and structure . . . . .	6
2.2 Data partitioning . . . . .	6
2.3 Machine-Learning approaches !WIP! . . . . .	6
2.3.1 K nearest neighbor . . . . .	7
2.3.2 Random forest . . . . .	7
2.3.3 Neural networks . . . . .	7
2.4 Quality metrics . . . . .	7
2.4.1 Terminology . . . . .	7
2.4.2 Visual metrics . . . . .	7
2.4.3 Accuracy . . . . .	7
2.4.4 False positive Rate . . . . .	7
2.4.5 Area under the curve . . . . .	8
2.4.6 Yield of Actives . . . . .	8
2.4.7 Enrichment Factor . . . . .	8
2.4.8 Relative Enrichment Factor . . . . .	8
2.5 Feature engineering !WIP! . . . . .	8
2.5.1 Feature engineering using random forest . . . . .	8
2.5.2 Physical properties . . . . .	8
2.5.3 Principal component analysis . . . . .	8
2.5.4 Balancing classes . . . . .	8

<b>3 Results !WIP!</b>	<b>9</b>
3.1 Performance per Protein-Complex . . . . .	9
3.2 Performance Overview – Comparing ML-approaches . . . . .	9
<b>4 Discussion</b>	<b>10</b>
4.1 Conclusion . . . . .	10
4.2 Improvements and outlook . . . . .	10
<b>A Technical Details</b>	<b>11</b>
<b>B Supplementary Materials</b>	<b>12</b>
B.1 PDF Files . . . . .	12
B.2 Media Files . . . . .	12
B.3 Online Sources (PDF Captures) . . . . .	12
<b>C Questionnaire</b>	<b>13</b>
<b>D LaTeX Source Code</b>	<b>14</b>
<b>References</b>	<b>15</b>
Literature . . . . .	15
Online sources . . . . .	17

# 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[1]. 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 generally be categorized as one of two methods. These two are called High-throughput-screening (HTS) and virtual screening (VS)[1].

When using an HTS-approach there are many compounds which are tested against some type of target protein. Target proteins are usually proteins which are of general interest for medical use. 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[3].

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[4]. 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[5].

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[6]. 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[4]. 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[5]. With flexible docking it is possible to simulate moving binding sites, where 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 conformation and 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[5].

Regardless of the docking type the score should reflect 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[7].

The focus of this work is on implementing a machine-learning based scoring approach. Machine-learning based scoring functions work by training on labeled 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[8]. In his thesis a selection of eleven proteins from the directory of useful decoys[9] have been selected to be analyzed. For the selected proteins all possible interactions have been analyzed by *PLIP*, which is an algorithm designed to discover various interactions based on the physical properties of a compound[10]. 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 drug discovery the state of the art is described in the following.

## 1.1 Machine Learning in drug design and activity prediction !WIP!

The following chapter summarizes the recent developments in drug design using various machine learning techniques.

Today there exist a multitude of machine learning approaches in the field of drug design and activity prediction. As a result of various AI breakthroughs in recent years there have been numerous research projects regarding the usability of artificial intelligence in various bioinformatic domains. MILCDock uses the Output of five traditional Scoring Functions as input for a neural Network. This technique has a slight performance benefit when compared to traditional scoring functions [11].

## 1.2 Goals

The goals of this thesis are twofold:

1. Evaluate common machine learning approaches for activity prediction and compare results with current literature.
2. Evaluate the results posed by various feature engineering techniques and investigate the possible performance benefits for the implemented ML approaches.

The second goal can be viewed as an extension of the first one since its primary aim is to improve the results achieved while pursuing the first goal.

## Chapter 2

# Methods

### 2.1 Data description

The following chapter is dedicated to explaining the data used for this thesis. This includes detailed descriptions of the protein complexes as well as their interaction types.

#### 2.1.1 Proteins

The following five proteins have been used to conclude this thesis:

##### Acetylcholinesterase

Acetylcholinesterase (AChE) is an efficient enzyme in the nervous system that breaks down acetylcholine (Ach), a messaging molecule, into choline and acetate. It's found in high concentrations at junctions between nerve cells and muscles. AChE has various functions beyond just breaking down Ach, and it's present in both nerve and non-nerve tissues. Because AChE is so important, some toxins like insecticides and nerve agents target it. This versatility of AChE makes it a key player in nervous system function and a potential target for drugs to treat diseases[12].

##### Cyclooxygenase 1

Cyclooxygenase 1 (COX1) and its isoform Cyclooxygenase 2 (COX2) play a substantial role in synthesising various prostaglandins. Due to their linkage with inflammations and pain COX molecules are often targeted by anti-inflammatory drugs. In contrast to COX2, COX1 is found in most tissues across the body. In addition to that COX1 is largely attributed with homeostatic functions such as hemostasis and gastric cytoprotection[13].

##### Dipeptidyl peptidase IV

Dipeptidyl peptidase IV(DPP4) can be partially responsible for hydrolysis of a prolyl bond between two residues from the N-terminus. DPP4 is present in several processes including metabolism and cancer biology. Due to its role within the metabolism DPP4 inhibitory drugs have been successfully used in the treatment of diseases type two. DPP4

also plays a substantial role in the diagnosis of certain types of cancer. In most cases DPP4 is up regulated near cancerous growth, therefore locally elevated DPP4 levels can be an indicator for cancer[14].

### Monoamine oxidase B

Monoamine oxidase B(MAOB) plays a major role in the breakdown of neurotransmitters(monoamines) within the body. The compound is mainly expressed in glial-cells and platelets. Its function categorizes MAOB as an important research compound, as MAOB inhibition has been proven to improve various neurological conditions. This stems from the fact that changes in the monoamine levels are associated with a myriad of neurological problems[15].

### Soluble epoxide hydrolase

Soluble epoxide hydrolase (sEH) is part of an inflammatory pathway similar to COX. It has been shown that inhibition of sEH reduces inflammation. In contrast to COX it does not completely disable the synthesis of pro-inflammatory compounds but rather balance their levels[16].

#### 2.1.2 Interactions → !WIP! ask Micha regarding grade of detail

The following interactions have been chosen for this thesis:

interaction	description[8]
hydrogen bonds	A hydrogen bond is defined as the interaction between a hydrogen atom, connected to a more electronegative atom, and another atom or molecule.
water bridges	A water bridge occurs when the ligand and the protein both bind to a water molecule through hydrogen bonds.
salt bridges	Salt bridges are ion pairs which stick together due to large difference in charge and the resulting electrostatic interaction.
halogen bonds	Halogen bonds are defined as the interactions between the electrophilic region around a halogen atom and a nucleophilic region.
hydrophobic interactions	Aggregates formed as a result of a hydrophobic interaction between hydrocarbons in an aqueous medium are called hydrophobic interactions.
pi-stacking	Interactions between neighboring aromatic rings are called pi-stacking. Due to the pi-electron density the ring is partially positively charged around the periphery and negatively charged above both aromatic faces. As a result electrostatic forces build between aromatic rings, and they are attracted to one another.
pi-cation	Cations and pi-stacks who bind through electrostatic forces at a pi-stacks face are called pi-cation interactions.

**Table 2.1:** interaction types

### 2.1.3 Data origin and structure

The provided data is a byproduct of the thesis [8] by Micha Birklbauer. The interaction data was produced using the *PLIP Algorithm* [10] on the aforementioned proteins.

The PLIP Algorithm consists of four major stages:

#### **Structural Preparation – SP**

During the preparation step the input structure is hydrogenated and the ligands(including their binding sites) are extracted.

#### **Functional Characterization – FC**

Using the structure of the complex a myriad of functional groups are detected. This includes binding site atoms, hydrophobic atoms and aromatic rings just to name a few.

#### **Rule Based Matching – RBM**

In the third step the algorithm investigates all interactions between the ligand and the protein, which can be attributed to geometric constraints. Hydrogen bonds are detected here.

#### **Filtering of Interactions – FoI**

This is a cleanup step where redundant or overlapping interactions get removed from the dataset.

The result of the PLIP Algorithm is a lineup of every interaction for each binding site and ligand[10]. This data has been used as a basis for the machine learning approaches discussed in this thesis.

## 2.2 Data partitioning

To validate the results of training various machine learning methods the provided data concerning the five targets was split into a training-set as well as a test-set. To achieve a 70/30 train/test split ratio each sample was randomly assigned to one of the two datasets[17].

In order to validate the machine learning approaches during training 10-fold cross-validation has been applied. For the process of cross-validation the training dataset is split into  $n$  equally large subsets. The type of cross-validation implemented in this thesis uses all but one of these partitions to train the classification model and validates the results with the remaining partition. This process is repeated for all possible validation partitions[18].

## 2.3 Machine-Learning approaches !WIP!

Introduction to the ML Approaches used for the thesis.

### 2.3.1 K nearest neighbor

### 2.3.2 Random forest

### 2.3.3 Neural networks

## 2.4 Quality metrics

In order to make the results from this thesis comparable to the results from the scoring function introduced in [8] by Micha Birklbauer the same quality metrics have been implemented for this thesis. The following will provide an overview for the used metrics.

### 2.4.1 Terminology

To calculate the metrics that are mentioned within this chapter a few base numbers are necessary:

**TP** – **T** rue **P** ositives are active samples, which are classified as such

**TN** – **T** rue **N** egatives are inactive samples, which are classified as such

**FP** – **F** alse **P** ositives are inactive samples, which are classified as active

**FN** – **F** alse **N** egatives are active samples, which are classified as inactive

### 2.4.2 Visual metrics

For better visualization of the four base metrics mentioned in 2.4.1 this thesis displays the resulting data in a confusion matrix. This metric displays distribution of the results over the four base metrics.

In addition to that, the ROC(receiver operating characteristic) curve will also be displayed for the results. The ROC curve is a collection of points in a two-dimensional space, where their location is defined by the FPR 2.4.4 on the x-axis and the TPR( $\frac{\#TP}{\#TP+\#FN}$ ) on the y-axis. Each point on this line depicts the ratio of FPR to TPR at a certain cutoff. [19]

### 2.4.3 Accuracy

Accuracy (ACC) describes which portion of the predicted samples was accurate and is defined as follows:

$$ACC = \frac{\#TP + \#TN}{\#TP + \#TN + \#FP + \#FN}$$

[20]

### 2.4.4 False positive Rate

False positive rate (FPR) describes the incorrect as active identified compounds in relation to all inactive compounds and is defined as follows:

$$FPR = \frac{\#FP}{\#TN + \#FP}$$

[19]



### 2.4.5 Area under the curve

Area under the curve (AUC) is a metric which stems from the ROC curve. The area under the ROC curve is calculated using the `scikit-learn` package.

### 2.4.6 Yield of Actives

Yield of actives (Ya) describes the correct as active identified compounds in relation to all as active labeled compounds and is defined as follows:

$$Y_a = \frac{\#TP}{\#TP + \#FP}$$

[21]

### 2.4.7 Enrichment Factor

The enrichment factor (EF) describes the relation of the truly active compounds among all as active predicted complexes and the relative share of active compounds in the dataset. This metric is defined as follows:

$$EF = \frac{\frac{\#TP}{\#TP + \#FP}}{\frac{\#TP + \#FN}{\#TP + \#TN + \#FP + \#FN}}$$

[19]

### 2.4.8 Relative Enrichment Factor

The relative enrichment factor (REF) describes the relation of the EF to the maximum achievable EF. The REF is defined as follows:

$$REF = \frac{100 * \#TP}{\min(\#TP + \#FP, \#TP + \#FN)}$$

[19]

## 2.5 Feature engineering !WIP!

Explain concept of feature engineering and possible implications for thesis.(Provide Overview here)

### 2.5.1 Feature engineering using random forest

### 2.5.2 Physical properties

### 2.5.3 Principal component analysis

### 2.5.4 Balancing classes

SMOTE

## Chapter 3

# Results !WIP!

### 3.1 Performance per Protein-Complex

Evaluate models tuned for selected proteins.

### 3.2 Performance Overview – Comparing ML-approaches

Compare performance of overall ml approaches.

## Chapter 4

# Discussion

### 4.1 Conclusion

Recap findings of thesis.

### 4.2 Improvements and outlook

Explain possible improvements to used technique. Provide general outlook on topic.

## 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 [22]

Appendix C

Questionnaire

Appendix D

LaTeX Source Code

# References

## Literature

- [1] 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).
- [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] 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).
- [4] 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 pp. 1, 2).
- [5] 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 pp. 1, 2).
- [6] 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).
- [7] 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).
- [8] Micha Johannes Birklbauer. “Automatic identification of important interaction-sand interaction-frequency-based scoring inprotein-ligand complexes”. MA thesis. FH Hagenberg, Aug. 31, 2021 (cit. on pp. 2, 5–7).
- [9] Michael M. Mysinger et al. “Directory of Useful Decoys, Enhanced (DUD-E): Better Ligands and Decoys for Better Benchmarking”. *Journal of Medicinal Chemistry* 55.14 (July 2012), pp. 6582–6594. DOI: 10.1021/jm300687e. (Visited on 04/30/2024) (cit. on p. 2).



- [10] Sebastian Salentin et al. “PLIP: fully automated protein–ligand interaction profiler”. *Nucleic Acids Research* 43.Web Server issue (July 2015), W443–W447. DOI: 10.1093/nar/gkv315. (Visited on 04/23/2024) (cit. on pp. 2, 6).
- [11] Connor Morris et al. “MILCDock: Machine Learning Enhanced Consensus Docking for Virtual Screening in Drug Discovery”. *Journal of chemical information and modeling* 62 (Nov. 2022). DOI: 10.1021/acs.jcim.2c00705 (cit. on p. 3).
- [12] Anurag Tripathi and U. C. Srivastava. “Acetylcholinesterase :A Versatile Enzyme of Nervous System”. *Annals of Neurosciences* 15.4 (Feb. 2010), pp. 106–111. DOI: 10.5214/95. (Visited on 04/22/2024) (cit. on p. 4).
- [13] Carol A. Rouzer and Lawrence J. Marnett. “Cyclooxygenases: structural and functional insights”. *Journal of Lipid Research* 50.Suppl (Apr. 2009), S29–S34. DOI: 10.1194/jlr.R800042-JLR200. (Visited on 04/22/2024) (cit. on p. 4).
- [14] Denise M. T. Yu et al. “The dipeptidyl peptidase IV family in cancer and cell biology”. en. *The FEBS Journal* 277.5 (2010), pp. 1126–1144. DOI: 10.1111/j.1742-4658.2009.07526.x. (Visited on 04/22/2024) (cit. on p. 5).
- [15] Rona R. Ramsay. “Molecular aspects of monoamine oxidase B”. *Progress in Neuro-Psychopharmacology and Biological Psychiatry* 69 (Aug. 2016), pp. 81–89. DOI: 10.1016/j.pnpbp.2016.02.005. (Visited on 04/22/2024) (cit. on p. 5).
- [16] Kara R. Schmelzer et al. “Soluble epoxide hydrolase is a therapeutic target for acute inflammation”. eng. *Proceedings of the National Academy of Sciences of the United States of America* 102.28 (July 2005), pp. 9772–9777. DOI: 10.1073/pnas.0503279102 (cit. on p. 5).
- [17] Yun Xu and Royston Goodacre. “On Splitting Training and Validation Set: A Comparative Study of Cross-Validation, Bootstrap and Systematic Sampling for Estimating the Generalization Performance of Supervised Learning”. en. *Journal of Analysis and Testing* 2.3 (July 2018), pp. 249–262. DOI: 10.1007/s41664-018-0068-2. (Visited on 04/21/2024) (cit. on p. 6).
- [18] Annette M. Molinaro, Richard Simon, and Ruth M. Pfeiffer. “Prediction error estimation: a comparison of resampling methods”. *Bioinformatics* 21.15 (May 2005), pp. 3301–3307. DOI: 10.1093/bioinformatics/bti499. eprint: [https://academic.oup.com/bioinformatics/article-pdf/21/15/3301/50340684/bioinformatics\\\_21\\\_15\\\_3301.pdf](https://academic.oup.com/bioinformatics/article-pdf/21/15/3301/50340684/bioinformatics\_21\_15\_3301.pdf) (cit. on p. 6).
- [19] Julio Cesar Dias Lopes et al. “The power metric: a new statistically robust enrichment-type metric for virtual screening applications with early recovery capability”. *Journal of Cheminformatics* 9.1 (Feb. 2017), p. 7. DOI: 10.1186/s13321-016-0189-4. (Visited on 04/20/2024) (cit. on pp. 7, 8).
- [20] Mohammad Hossin and Sulaiman M.N. “A Review on Evaluation Metrics for Data Classification Evaluations”. *International Journal of Data Mining & Knowledge Management Process* 5 (Mar. 2015), pp. 01–11. DOI: 10.5121/ijdkp.2015.5201 (cit. on p. 7).
- [21] Deborah Giordano et al. “Drug Design by Pharmacophore and Virtual Screening Approach”. *Pharmaceuticals* 15.5 (May 2022), p. 646. DOI: 10.3390/ph15050646. (Visited on 04/20/2024) (cit. on p. 8).

## Online sources

- [22] *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! —