

MASARYKOVA
UNIVERZITA

MASARYK UNIVERSITY
Faculty of Science



Ph.D. Thesis

Title -

Brno 2021

Stanislav Geidl

To my grandma.

Acknowledgements

Declaration

I hereby declare that this disertation thesis is my original authorial work, which I have worked on alone. All sources, references and literature used or excerpted during the elaboration of this work are properly cited and listed in a complete reference with regard to the source.

Brno, xxth xxxxxxxx 2021 Stanislav Geidl

Contents

I	Introduction	1
1	Introduction	2
II	Theory and Methods	4
2	Structure	5
2.1	Molecular Structure in Computer	5
2.2	3D Structure Calculation	5
2.2.1	Rule-Based and Data-Based Methods	5
2.2.2	Fragment-Based Method	5
2.2.3	Numerical Method	5
2.2.4	Conformational Analysis	5
3	Partial Atomic Charges	6
3.1	The Concept of Atomic Charges	6
3.2	Overview of Charge Calculation Methods	6
3.2.1	QM Charge Methods	6
3.2.2	Empirical Methods	6
3.3	EEM Calculation	6
3.4	Quality and Usability of EEM parameters	6
3.5	EEM Parametrization	6
4	Acid Dissociation Constant Prediction	7
4.1	Motivation	7
4.2	Overview of Approaches	7
4.2.1	LFER (Linear Free Energy Relationships) Methods	7
4.2.2	Database Methods	7
4.2.3	Ab Initio Quantum Mechanical Calculations	7
4.2.4	QSPR Method	7
III	Results	8
5	Synopsis of the Results	9

6	Follow-up work and future plans	10
IV	Conclusion	11
7	Conclusion	12
	Bibliography	14
	Appendix: Main papers	15
	Appendix: All My Publications	16
	Curriculum Vitae	17

Part I

Introduction

Introduction

In recent years, a vast amount of data about various types of molecules became available. For example, we can obtain the complete human genome of a selected individual in a few days, and about 150 thousand biomacromolecular structures have been determined and published (Protein Data Bank [1]). Furthermore, more than 100 million various small molecules are described in freely accessible databases (e.g., Pubchem [2], ZINC [3], ChEMBL [4]). This richness of data caused the formation of novel modern life-science research fields focused on the utilization of this data. The best-known modern life sciences are bioinformatics, structural bioinformatics, systems biology, genomics, proteomics, and also chemoinformatics. These current research specializations have provided many key results in basic and applied research (e.g. [5–11]).

One fascinating and beneficial field utilizing and processing newly available data about small molecules (i.e., drug-like compounds) is chemoinformatics. This discipline offers methodologies for comparing molecular similarity, molecular database search, virtual screening, and the prediction of molecules' properties and activities. This prediction is based on the idea that molecular structures' similarity has a consequence – a similarity in molecular properties. In chemoinformatics, the structure is first described using mathematical characteristics (so-called descriptors) – numbers containing 3D (or 2D or 1D) structure information and applicable as inputs of mathematical models. Then, these models are constructed based on a relation between descriptors and known values of the property or the activity. Such models are called Quantitative Structure-Property Relationship (QSPR) models or Quantitative Structure-Activity Relationship (QSAR) models.

A property, which is strongly required and is therefore often a target of chemoinformatics prediction models is the acid dissociation constant, K_a , and its negative logarithm pK_a . Those pK_a values are of interest in chemical, biological, environmental, and pharmaceutical research [12–14]. pK_a values have found applications in many areas, such as evaluating and optimizing drug candidate molecules, pharmacokinetics, ADME profiling, understanding protein-ligand interactions, etc. Moreover, the critical physicochemical properties such as permeability, lipophilicity, solubility, etc., are pK_a dependent. Unfortunately, experimental pK_a values are available only for a limited set of molecules. In addition to that, obtaining experimental pK_a values for newly designed molecules is very time-consuming because they must be synthesized first. Chemoinformatics approaches for pK_a prediction are therefore currently intensively examined.

For this reason, I also focused on the chemoinformatics way of pK_a prediction in my work. Very promising descriptors for pK_a prediction are partial atomic charges [] because they hold

information about the distribution of electron density within the molecule. Specifically, electron densities on atoms close to the dissociating hydrogen provide a clue about its dissociation ability. The most common and accurate method for calculating partial atomic charges is an application of quantum mechanics (QM). QM calculation can be performed via various approaches, introducing different approximation levels (i.e., approximating a wave function by different sets of mathematical equations, which are called basis sets). QM outputs electron distribution in orbitals and this distribution can be divided into individual atoms using several charge calculation schemes (e.g., MPA, NPA, AIM, Hirshfeld, MK, etc.). Therefore, the correlation between pK_a and relevant atomic charges calculated by different QM approaches has been analyzed []. I also focused on this file in my bachelor thesis, developed a workflow for calculation of pK_a using QM partial atomic charges and examined, which types of QM are the most suitable [].

QM charges are accurate, but their calculation is very time-consuming. A faster Alternative to QM charges is empirical charge calculation approaches. Furthermore, if we would like to apply chemoinformatics pK_a prediction models practically – for example, in pre-screening large sets of drug candidates – we need a fast approach. Therefore, in my master thesis, I developed a pK_a prediction workflow based on charges (including Electronegativity Equalization Method).

However, several pieces of the puzzle were still missing. For example, the developed pK_a prediction workflows [] were strongly dependent on 3D structure source, and also, the quality of available EEM charges was low.

Therefore, my dissertation’s goal was to develop a workflow that predicts pK_a for molecules not synthesized yet and without available experimental 3D structures.

Specifically, the thesis examined how to improve the process of pK_a prediction via providing suitable inputs. First, the influence of 3D structure source on pK_a prediction accuracy was analyzed. Afterward, the work focused on obtaining high-quality partial atomic charges, which served as descriptors for pK_a calculation. In the end, the authors also support the development of methodology and software tools for obtaining these high-quality charges.

The thesis structure is the following: First, an overview of key fields is provided (Part II), i.e. – 3D structure and approaches for its prediction, charge calculation methods, and pK_a prediction approaches. Next, the achieved results, which we published in three research papers, are briefly described (Part III), and full-texts of the respective published papers are attached in Appendix: Main papers. During the elaboration of this thesis, I was also involved in other projects. Most of them were not related to pK_a prediction but tightly connected to the field of chemoinformatics or structural bioinformatics. The outcome of these projects consists of several papers and a book I have co-authored. Their title pages are attached in Appendix: All My Publications.

Part II

Theory and Methods

2

Structure

2.1 Molecular Structure in Computer

2.2 3D Structure Calculation

2.2.1 Rule-Based and Data-Based Methods

2.2.2 Fragment-Based Method

2.2.3 Numerical Method

2.2.4 Conformational Analysis

Partial Atomic Charges

- 3.1 The Concept of Atomic Charges**
- 3.2 Overview of Charge Calculation Methods**
 - 3.2.1 QM Charge Methods**
 - 3.2.2 Empirical Methods**
- 3.3 EEM Calculation**
- 3.4 Quality and Usability of EEM parameters**
- 3.5 EEM Parametrization**

Acid Dissociation COnstant Prediction

4.1 Motivation

4.2 Overview of Approaches

4.2.1 LFER (Linear Free Energy Relationships) Methods

4.2.2 Database Methods

4.2.3 Ab Initio Quantum Mechanical Calculations

4.2.4 QSPR Method

Part III

Results

5

Synopsis of the Results

6

Follow-up work and future plans

Part IV

Conclusion

7

Conclusion

Appendix

Bibliography

- [1] Berman, H. M., Kleywegt, G. J., Nakamura, H., and Markley, J. L. (2014) The protein data bank archive as an open data resource. *Journal of computer-aided molecular design*, **28**, 1009–1014.
- [2] Bolton, E. E., Wang, Y., Thiessen, P. A., and Bryant, S. H. (2008) Pubchem: integrated platform of small molecules and biological activities. *Annual reports in computational chemistry*, **4**, 217–241.
- [3] Irwin, J. J., Sterling, T., Mysinger, M. M., Bolstad, E. S., and Coleman, R. G. (2012) Zinc: A free tool to discover chemistry for biology. *Journal of chemical information and modeling*, **52**, 1757–1768, PMID: 22587354.
- [4] Gaulton, A., et al. (2012) ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic acids research*, **40**, D1100–D1107.
- [5] Paulsen, C. E., Armache, J.-P., Gao, Y., Cheng, Y., and Julius, D. (2015) Structure of the TRPA1 ion channel suggests regulatory mechanisms. *Nature*, **520**, 511–517.
- [6] Cao, E., Liao, M., Cheng, Y., and Julius, D. (2013) Trpv1 structures in distinct conformations reveal activation mechanisms. *Nature*, **504**, 113–118.
- [7] Protá, A. E., Bargsten, K., Zurwerra, D., Field, J. J., Díaz, J. F., Altmann, K.-H., and Steinmetz, M. O. (2013) Molecular mechanism of action of microtubule-stabilizing anticancer agents. *Science*, **339**, 587–590.
- [8] Lu, J., et al. (2005) MicroRNA expression profiles classify human cancers. *Nature*, **435**, 834–838.
- [9] Puente, X. S., et al. (2011) Whole-genome sequencing identifies recurrent mutations in chronic lymphocytic leukaemia. *Nature*, **475**, 101–105.
- [10] Nayal, M. and Honig, B. (2006) On the nature of cavities on protein surfaces: application to the identification of drug-binding sites. *Proteins*, **63**, 892–906.
- [11] Xie, L., Xie, L., and Bourne, P. (2009) A unified statistical model to support local sequence order independent similarity searching for ligand-binding sites and its application to genome-based drug discovery. *Bioinformatics*, **25**, i305–i312.
- [12] Comer, J. and Tam, K. (2001) *Pharmacokinetic Optimization in Drug Research: Biological, Physicochemical, and Computational Strategies*. Verlag Helvetica Chimica Acta, Postfach, CH-8042 Zürich, Switzerland.
- [13] Klebe, G. (2000) Recent developments in structure-based drug design. *Journal of molecular medicine*, **78**, 269–281.
- [14] Kim, J. H., Gramatica, P., Kim, M. G., Kim, D., and Tratnyek, P. G. (2007) Qsar modelling of water quality indices of alkylphenol pollutants. *SAR and QSAR in environmental research*, **18**, 729–743.

Appendix: Main papers

Appendix: All My Publications

Curriculum Vitae