

Charles University
Faculty of Science

Study programme: Bioinformatics

Branch of study: Bioinformatics



Kateřina Břicháčková

Use of residue-level annotations for structural
prediction of protein-ligand binding sites

Využití anotací primární struktury pro strukturní
predikci protein-ligand aktivních míst

Master thesis

Supervisor: RNDr. David Hoksza, Ph.D.

Prague, 2020

Prohlášení

Prohlašuji, že jsem závěrečnou práci zpracovala samostatně a že jsem uvedla všechny použité informační zdroje a literaturu. Tato práce ani její podstatná část nebyla předložena k získání jiného nebo stejného akademického titulu.

V Praze, XXX

Kateřina Břicháčková

Poděkování

Ondrovi :P

Acknowledgement

Ondrovi :P

Abstract

abstract

Keywords: keyword1, keyword2

Abstrakt

abstract

Klíčová slova: keyword1, keyword2

Contents

1	Introduction	3
2	Ligand binding sites predictors	4
3	Methods	5
3.1	Statistical analysis	5
3.1.1	Welch’s test	6
3.1.2	Fisher’s exact test and Chi-squared test of independence .	8
3.1.3	Implementation in Python	9
3.2	Pipeline description	9
4	Datasets	10
5	Features	11
5.1	UniProtKB	11
5.1.1	PTM	11
5.1.2	Disulfide bonds	12
5.1.3	Non-standard residues	12
5.1.4	Secondary structure	12
5.1.5	Natural variant	12
5.1.6	Variation	12
5.1.7	Compositional bias	12
5.2	PDBe-KB	12
5.2.1	Conservation	13
5.2.2	DynaMine	13
5.2.3	EFoldMine	13
5.2.4	Depth	14
5.3	PDB	14
5.3.1	B factor	14
5.3.2	Half sphere exposure	14
5.3.3	Exposure CN	14
5.3.4	Phi and psi angles	14
5.3.5	Cis peptide	14
5.4	FASTA	14
5.5	Other resources	14
5.5.1	MobiDB	14
5.5.2	Conservation	14
6	Results	15
	Conclusion	16
	List of Abbreviations	17
	Bibliography	18

A	Attachments	20
A.1	attachment1	20
A.2	attachment2	20

1. Introduction

2. Ligand binding sites predictors

...

3. Methods

TODO verze programu a databazi? Kam to dat?

To find properties mapped on the protein primary structure which are possibly important for prediction of protein-ligand binding sites, statistical analysis will have the crucial role. It is a great way to explore the big amounts of accessible data and it can potentially help to discover underlying patterns and draw inferences from the data.

This chapter describes the method that was used to analyse the *statistical significance* of the properties and to distinguish the ones that stand out in the known protein-ligand binding sites.

3.1 Statistical analysis

Hypothesis testing is a method of statistical inference. Its goal is to infer properties of a *statistical population*, i.e. a set of similar items or events. In this work, two populations will be compared: we take values of a property for all the amino acids across all the proteins in the dataset and then compare the ones in binding sites and outside of binding sites.

A dataset usually contains a subset sampled from a larger population, rather than the whole population. This subset is called a *statistical sample*. It should represent the population well and be unbiased.

A *hypothesis* makes a statement about an unknown population parameter. In a hypothesis testing problem, an experimenter states two complementary hypotheses: the *null hypothesis* and the *alternative hypothesis*, denoted by H_0 and H_1 , respectively. The null hypothesis comprises a subset of possible parameters and the alternative hypothesis comprises the supplement, so that all the possible parameters are covered.

In a hypothesis testing problem, an experimenter should come to one of the conclusions: to either accept H_0 , or to reject H_0 and accept H_1 .

To decide which one of two complementary hypotheses is true, an experimenter employs a suitable *hypothesis test*. A hypothesis test is a rule that specifies for which sample values the H_0 is accepted as true and for which sample values it is rejected, and therefore H_1 is accepted as true. A hypothesis test is usually specified in terms of a test statistic (i.e. a function of the sample) [3].

As one may expect, the tests are not error-proof and a mistake can be made in the decision of whether to accept or reject the null hypothesis. There are two types of errors in hypothesis testing, commonly known as *Type I error* and *Type II error*. The test has made a Type I error if it incorrectly rejects a true null hypothesis. If, on the other hand, a null hypothesis is accepted and it is not true, a Type II error has been made. Both situations are depicted in the Table 3.1. The ideal test would have both error probabilities equal to zero. Nevertheless, in most cases it is not possible to make both error probabilities arbitrarily small for a fixed sample size [3].

To control statistical significance of the result, a study defines a threshold called *significance level*, a constant denoted by α . It represents the probability of

		Prediction	
		Accept H_0	Reject H_0
Truth	H_0	Correct (true positive)	Type I error (false positive)
	H_1	Type II error (false negative)	Correct (true negative)

Table 3.1: Type I and II Error in hypothesis testing.

making a Type I error, in other words, the probability that the study rejects the null hypothesis when it is true.

One way to report the result of the test would be simply to tell whether the null hypothesis was accepted or rejected at the given significance level. However, most researchers choose to report a certain kind of test statistic (function of a sample X), the so-called *p-value*. Smaller values of $p(X)$ give stronger evidence for rejecting the null hypothesis. The null hypothesis is rejected when $p(X) \leq \alpha$. Hence, we are able to determine the smallest significance level at which the hypothesis would be accepted/rejected. P-value gives an idea of how strongly the data contradict the null hypothesis; furthermore, it allows other researchers to make a decision according to the significance level of their choice [3, 5, 11].

It is suggested that the significance level for a study is set prior to any data collection [9]. The typical choices in practice are $\alpha = 0.01, 0.05$ or 0.10 [3]. One should be aware that by fixing the significance level of the test, the experimenter is controlling only the Type I error probabilities. The probability of the Type II error is subject to factors such as the accuracy and completeness of the data and most importantly, the true effect size [11].

Let's suppose an experimenter has a research hypothesis that he or she hopes to prove, but does not want to risk accepting it without convincing data support. In this case, the test should be set up in such a way that the research hypothesis corresponds to the alternative hypothesis, not the null hypothesis. By specifying a small significance level α , the experimenter thus controls the probability of the Type I error. In other words, the probability of accepting the research hypothesis when it is not true would be α at most [3].

3.1.1 Welch's test

Welch's unequal variances t-test, or Welch's test in short, is a two-sample hypothesis test used to decide whether two populations have different central tendencies (means or medians). The decision is made based on the samples from the two populations. It is a more robust alteration of the widely used Student's t-test [12].

Both Student's and Welch's t-test assume that the two examined populations follow a normal distribution [12]. Nevertheless, when testing for the equality of means of "large enough samples", the normality assumption can be violated

thanks to the large sample theory and the Central Limit Theorem [5]. It has been shown in previous studies that for large samples, the statistical significance level is protected not only for normally distributed data, but also for many non-normal distributions; moreover, in case of Welch’s test, this is true even for unequal variances [7, 15, 16]. According to Lehmann and Romano [5], the Type II error is also relatively insensitive to non-normality. Many articles and textbooks mention that when the sample sizes are small, nonparametric tests (i.e. tests that do not assume a specific distribution) such as the Mann-Whitney test [?] should be considered as an alternative to t-tests. However, t-tests become superior when sample sizes increase [7, 13]. The simulations made by Lumley *et al.* [7] show that “sufficiently large sample size” means under 100 in most cases. Even for extremely non-normal data, the sufficient size is at most 500. This suggests that the choice of Welch’s test is legitimate for this work.

The problem of the Student’s t-test is that it performs badly when the variances of the two compared populations are unequal. Both Type I and Type II errors are negatively affected by violation of the equal variances assumption. The unequal variances can be less problematic if sample sizes are similar, but in practice, that is not always the case [10].

Unlike Student’s t-test, Welch’s test does not assume equal variances of the populations. It performs well when the samples have unequal variances; furthermore, it can be used even when the samples have unequal sizes [4].

Some researchers tend to pre-test for variance equality by a preliminary test of variances (such as Levene’s [6], Bartlett’s [?] or Brown-Forsythe test [?]) and then choose whether to use Student’s or Welch’s t-test. However, although this approach persists in some textbooks and software packages, it is not recommended by statisticians. As a preliminary test itself is subject to Type I and II errors, this two-stage procedure would not protect the significance level and could lead to incorrect decisions. One should be aware of the fact that even if the test suggested that the samples variances are nearly equal, it would not mean that the whole population variances could not differ to a larger extent [14]. Some researchers may try to make the significance level of a preliminary test more strict, so that they could be more confident about the choice of the subsequent test; however, as the significance level decreases, the performance of the compound test paradoxically gets worse. According to Zimmerman [14], “a higher Type I error rate of the preliminary test actually improves the performance of the compound test”. This suggests that using the preliminary test is not correct in principle.

Welch’s test should be used whenever the researcher is not sure that the variances are truly equal. Ruxton [10] even suggests the routine use of Welch’s test. When the sample sizes and variances are equal, both tests perform similarly. When dealing with unequal variances and unequal sample sizes, Welch’s test is more robust than Student’s t-test and the Type I error rate does not deviate far from the nominal value [4]. Hence, Welch’s test can be applied without any significant disadvantages to Student’s t-test.

For all the reasons stated above, Welch’s test seems to be the best choice for the purpose of this study. It has the best combination of performance and ease of use, the calculation is straightforward and it is available in commonly used statistics packages.

3.1.2 Fisher’s exact test and Chi-squared test of independence

A different kind of tests will be needed for the analysis of categorical features. An example of a categorical feature is XXX. Moreover, quantitative features can be grouped into categories and analysed in the same way as categorical features. In this section, two widely-used tests of such kind will be presented and discussed.

Both Fisher’s exact test and Chi-squared test of independence are well-known hypothesis tests used for the analysis of data in contingency tables. A *contingency table* is a table displayed in a form of a matrix where cells represent a frequency distribution of samples in the categories. An example of a contingency table can be seen in Table 3.2. The sums of frequencies in rows and columns are called *marginal totals*.

	PTM	Without PTM	Total
Binding sites	XX	XX	XX
Non-binding sites	XX	XX	XX
Total	XX	XX	XX

Table 3.2: A 2×2 contingency table. TODO real data

The null hypothesis assumes independence of the groups; in our case, the assumption is that there is no difference in the proportions of the analysed feature between binding sites and non-binding sites.

Fisher’s exact test belongs to a class of so-called *exact tests*; it means that the p-value is calculated accurately, not approximately, as is the case of many tests including Welch’s test and Chi-squared test. Fisher’s test is mostly used for 2×2 contingency tables, although the principle of the computation can be extended to a general $m \times n$ table [8]. The principle of the test lies in computing the probability of obtaining a table that is more or equally extreme in the departure from the null hypothesis than the analysed table and has identical marginal totals [2].

Chi-squared (χ^2) test if independence is able to decide whether the difference between the observed frequencies and the “expected frequencies” is statistically significant. The expected frequencies are computed for every cell as $\frac{\text{row total} \times \text{column total}}{\text{grand total}}$. It can be imagined as the average frequencies we would get in the long run with the same marginal totals, assuming the null hypothesis is true (i.e. there is no association between groups). The result of the test tells how likely are we to observe given data under the assumption of the true null hypothesis [2]. TODO vysvetlit chi-sq rozdeleni a to jak se pocita ta testova statistika

The biggest difference between the two mentioned tests is that the chi-squared test is based on a aproximation approach; therefore, it needs a “large enough” sample. TODO FISHERUV TEST OBEČNE NA MENŠI VZORKY, W. G. Cochran (1954) [?] proposed a set of recommendations about the minimum expectations to be used in χ^2 tests and about the choice between Fisher’s test and χ^2 test:

These recommendations are presented in several textbooks and articles as a rule of thumb [] and recommended to be used in practice.

TODO: for small, sparse, or unbalanced data, the exact and asymptotic p-values can be quite different and may lead to opposite conclusions concerning the hypothesis of interest. (wikipedia)

3.1.3 Implementation in Python

3.2 Pipeline description

4. Datasets

The choice of datasets of protein-ligand complexes used for statistical analysis and P2Rank model training and evaluation was strongly inspired by the datasets described in the P2Rank article [?]. The structures were re-downloaded directly from PDBe, according to their PDB ID (four-character alphanumeric identifier) and chain ID (one-character identifier) used in the original datasets. It was not possible to take the original datasets as they were, since the structures were not up-to-date and the annotations downloaded from the databases (e.g. feature values) could not be mapped properly.

Downloaded datasets were further checked and filtered: Obsolete structures were replaced with their current entries, structures that do not have a corresponding UniProt record were removed, as well as structures with the incorrect segments mapping due to the bug in PDBe (mentioned in section TODO-ODKAZ).

The resulting datasets were named identically with the original datasets:

- **Chen11** - a smaller non-redundant dataset that was originally designed for a comparative study of ligand binding sites predictors [?]. It comprises at most one representative chain for every SCOP family [?] to ensure the minimal sequence similarity and maximal variability in tertiary structure. The original dataset covers 6 structural classes, 148 protein folds, 184 superfamilies and 251 families [?]; after re-downloading and filtering, the numbers are slightly smaller.
- **Coach420** - a dataset that was originally taken from a benchmark study [?] and used in other studies [? ?]. The non-redundant dataset harbor mix of natural and drug-like ligand molecules.
- **Joined** - smaller datasets from previous studies merged together in one larger dataset. It comprises a set of drug-target complexes extracted from DrugBank, DrugPort and PDB DT198[?], a benchmark set for the validation of protein-ligand docking performance [?], and a dataset with bound and unbound structures used for evaluation of a ligand binding sites predictor [?].
- **Holo4k** - a large set of protein-ligand complexes used in a large-scale evaluation of four binding sites predictors [?].

TODO single chain

5. Features

....TODO

5.1 UniProtKB

The UniProt Knowledgebase (UniProtKB) is a large database of well-annotated protein sequence data. It tries to achieve the minimal redundancy and it provides detailed, accurate and consistent annotations of the sequences [?].

Sequence annotations (called ‘features’) are available for every UniProtKB entry. They describe interesting sites and regions on the protein sequence and every feature has an associated description with more information, such as available evidence, source or related publications. The features are arranged in a well-organized manner on the website, in so called ‘Features viewer’ with many overlapping tracks for different features. Nonetheless, for the purpose of this work, the best way to obtain the features was via the Proteins REST API [?]. It provides the interface to access the sequence annotation data as well as mapped variation data programmatically. The API is available at (<http://www.ebi.ac.uk/proteins/api/doc>) [?].

Features are classified into eight categories which are further subdivided into types. For example, the category ‘STRUCTURAL’ comprises the types ‘HELIX’, ‘TURN’ and ‘STRAND’.

The types and categories that were chosen as potentially relevant for ligand binding sites prediction are described below.

5.1.1 PTM

Post-translational modifications are covalent chemical modifications of polypeptide chains after translation, usually modifying the functional group of the standard amino acids, or introducing a new group. They extend the set of the 20 standard amino acids and they can be important for the function of many proteins, as they can alter the interactions with other proteins, localization, activity, signal transduction, cell-cell interactions and other properties. Their enrichment in binding sites is very interesting to examine.

Three UniProtKB feature types were analysed: lipidation, glycosylation and type ‘MOD.RES’ which comprises phosphorylation, methylation, acetylation, amidation, formation of pyrrolidone carboxylic acid, isomerization, hydroxylation, sulfation, flavin-binding, cysteine oxidation and nitrosylation. Only experimentally determined modification sites are annotated, and they are further propagated to related orthologs when specific criteria are met [?].

Since lipidation and glycosylation data were very sparse (e.g. there were only 15 lipidation sites in the whole hol4k dataset composed of 3973 proteins), the fourth feature called ‘PTM’ including all three types was added to the analysis.

5.1.2 Disulfide bonds

Another type of post-translational modifications are disulfide bonds formed between two cysteine residues. Both intrachain and interchain bonds are annotated by UniProtKB. The disulfide bonds may be either experimentally determined or predicted [?].

5.1.3 Non-standard residues

Describes the occurrence of non-standard amino acids (selenocysteine and pyrrolysine). There must be experimental evidence for this occurrence; however, it can be propagated to close homologs [?].

5.1.4 Secondary structure

This feature category annotates three types of secondary structures: helices, beta sheets and hydrogen-bonded turns. Residues not belonging to any of the classes are in a random-coil structure. The ‘helix’ class comprises alpha-helices, pi-helices and 3_{10} helices.

The secondary structure assignment is made by DSSP algorithm [?] based on the coordinate data sets extracted from the Protein Data Bank (PDB). They are neither predicted computationally, nor propagated to related species [?].

5.1.5 Natural variant

This feature includes naturally occurring polymorphisms, variations between strains or RNA editing events [?].

5.1.6 Variation

Variation service is a utility that can retrieve variation data from UniProtKB. The variants are either extracted from the scientific literature and manually reviewed, or mapped from large scale studies, such as 1000 Genomes [?], COSMIC [?], ClinVar [?] or ExAC [?]. The Proteins REST API provides various options for variants retrieval, such as to filter by the consequence type, associated disease name, cross reference database type (e.g. ClinVar) or by the source type [?].

5.1.7 Compositional bias

The regions of compositional bias are parts of the polypeptide chain where some of the amino acids are over-represented, not following the standard frequencies. The regions can be enriched in one or more different amino acids [?].

5.2 PDBe-KB

PDBe-KB (Protein Data Bank in Europe - Knowledge Base) is managed by the PDBe team at the European Bioinformatics Institute. It is a collaborative resource that aims to bring together the annotations from various sources and to show the macromolecular structures in broader biological context.

One drawback of PDB is that every page represents only one entry that is based on a single experiment. There may be several PDB entries for the full-length protein, each covering only a segment of it. Nevertheless, the entries for the same protein are not interconnected. PDBe-KB has developed the *aggregated views of proteins*, displaying an overview of all the data related to the full-length protein defined by the UniProtKB accession.

The structures from the PDB are extensively used by scientific software and other resources. There exist many valuable annotations, such as ligand binding sites, post-translational modification sites, molecular channels or effects of mutations, that are created outside of the PDB. The problem is that the data is fragmented and therefore it would require immense effort of a researcher to collect and make use of all available data for a structure of interest.

The aggregated views of proteins integrates the annotations from *PDBe-KB partners*, collaborating scientific software developers. It facilitates the retrieval of these annotations with a uniform data access mechanism (via FTP or REST API). The project is called ‘FunPDBe’. A common data exchange scheme was defined to facilitate the transfer of data. [?]

The use of PDBe-KB was difficult because of the lack of documentation and a few bugs that were encountered during this work (some of them corrected by now after pointing them out). However, it is understandable since it was launched only two years ago and the constant improvements are done since then.

5.2.1 Conservation

PDBe-KB provides pre-calculated residue-level conservation scores, obtained by a pipeline using HMMER and Skyline web servers that was described by Jakubec et al. [?].

The values of the score are integers ranging from 0 to 9, with 9 being the most conserved. Since scores higher than 4 were very sparse and the feature would not meet the assumptions of the Chi-squared test, the scores 4 and higher were merged into one category (4). This does not deteriorate the prediction nor the hypothesis test, as vast majority (over 95%) of non-binding residues were scored 1 and lower.

5.2.2 DynaMine

DynaMine was developed by the Bio2Byte group [?] and it is one of the PDBe-KB partner resources. It provides the annotations of the backbone dynamics predicted only from the FASTA sequence. DynaMine predicts backbone flexibility at the residue-level, using a linear regression model trained on a large dataset of curated NMR chemical shifts extracted from the Biological Magnetic Resonance Data Bank [?]. The predictor estimates the value of the ‘order parameter’ (S^2) which is related to the rotational freedom of the N-H bond vector of the backbone. The values range from 0 (highly dynamic) to 1 (complete order) [?].

5.2.3 EFoldMine

EFoldMine tool comes from the same group as DynaMine. It is a predictor of the early folding regions of proteins. It makes predictions at the residue-level

derived only from the FASTA sequence. Internally it uses dynamics predictions and secondary structure propensities as features and the linear regression model is trained on data from NMR pulsed labelling experiments. Unfortunately, the early stages of protein folding are not understood very well so far and experimental data is very difficult to obtain. The predictor was trained on the dataset of only 30 proteins and its performance is quite poor [?].

5.2.4 Depth

Depth is a webserver that can measure residue burial within the protein. It is able to find small cavities in proteins and could be used as a ligand-binding sites predictor as such. The residue depth values are computed from the input PDB file.

The algorithm places input 3D structure in the box of model water, each residue with at least two hydration shells around itself. The water molecules in cavities are removed: the algorithm removes the water molecule if there are less than a given number of water molecules in its spherical volume of given size. The minimum number of neighbouring molecules and the spherical volume can be defined by the user. The removal is iterated until there are no more cavity waters. Residue depth is then computed as the distance to the closest water molecule [?].

5.3 PDB

5.3.1 B factor

5.3.2 Half sphere exposure

5.3.3 Exposure CN

5.3.4 Phi and psi angles

5.3.5 Cis peptide

5.4 FASTA

5.5 Other resources

5.5.1 MobiDB

5.5.2 Conservation

6. Results

...

Conclusion

Conclusion.

List of Abbreviations

AA Amino acid
atd a tak dale

Bibliography

- [1]
- [2] Martin Bland. *An introduction to medical statistics*. Oxford University Press, Oxford New York, 1987. ISBN 0192615025.
- [3] George Casella and Roger L Berger. *Statistical inference*, volume 2. Duxbury Pacific Grove, CA, 2002.
- [4] B. Derrick and P. White. Why welch's test is type i error robust. *The Quantitative Methods for Psychology*, 12(1):30–38, jan 2016. doi: 10.20982/tqmp.12.1.p030.
- [5] Joseph P. Romano Erich L. Lehmann. *Testing Statistical Hypotheses*. Springer New York, 2008. ISBN 0387988645. URL https://www.ebook.de/de/product/3186913/erich_l_lehmann_joseph_p_romano_testing_statistical_hypotheses.html.
- [6] Howard Levene. Robust tests for equality of variances. In Ingram Olkin, editor, *Contributions to Probability and Statistics: Essays in Honor of Harold Hotelling*, pages 278–292. Stanford University Press.
- [7] Thomas Lumley, Paula Diehr, Scott Emerson, and Lu Chen. The importance of the normality assumption in large public health data sets. *Annual Review of Public Health*, 23(1):151–169, may 2002. doi: 10.1146/annurev.publhealth.23.100901.140546.
- [8] Cyrus R. Mehta and Nitin R. Patel. A network algorithm for performing fisher's exact test in $r \times c$ contingency tables. *Journal of the American Statistical Association*, 78(382):427, jun 1983. doi: 10.2307/2288652.
- [9] Jerzy Neyman and Egon S Pearson. *The testing of statistical hypotheses in relation to probabilities a priori*, volume 29. 1933.
- [10] Graeme D. Ruxton. The unequal variance t-test is an underused alternative to student's t-test and the mann–whitney u test. *Behavioral Ecology*, 17(4): 688–690, may 2006. doi: 10.1093/beheco/ark016.
- [11] Pak C. Sham and Shaun M. Purcell. Statistical power and significance testing in large-scale genetic studies. *Nature Reviews Genetics*, 15(5):335–346, apr 2014. doi: 10.1038/nrg3706.
- [12] B. L. WELCH. THE GENERALIZATION OF ‘STUDENT’S’ PROBLEM WHEN SEVERAL DIFFERENT POPULATION VARLANCES ARE INVOLVED. *Biometrika*, 34(1-2):28–35, 1947. doi: 10.1093/biomet/34.1-2.28.
- [13] Donald W. Zimmerman. Invalidation of parametric and nonparametric statistical tests by concurrent violation of two assumptions. *The Journal of Experimental Education*, 67(1):55–68, jan 1998. doi: 10.1080/00220979809598344.

- [14] Donald W. Zimmerman. A note on preliminary tests of equality of variances. *British Journal of Mathematical and Statistical Psychology*, 57(1):173–181, may 2004. doi: 10.1348/000711004849222.
- [15] Donald W. Zimmerman and Bruno D. Zumbo. Rank transformations and the power of the student t test and welch t' test for non-normal populations with unequal variances. *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie expérimentale*, 47(3):523–539, 1993. doi: 10.1037/h0078850.
- [16] Bruno D. Zumbo and Daniel Coulombe. Investigation of the robust rank-order test for non-normal populations with unequal variances: The case of reaction time. *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie expérimentale*, 51(2):139–150, 1997. doi: 10.1037/1196-1961.51.2.139.

A. Attachments

The attached CD contains two attachments:

A.1 attachment1

blabla

A.2 attachment2

blabla