



VIREX



IŞIK UNIVERSITY
COMPUTER
SCIENCE AND
ENGINEERING

Gözde Gökyokuş Supervised by Asst. Prof. İlknur Karadeniz

Introduction

As the world is facing a dangerous pandemic nowadays, we clearly realize the importance of medical services and the importance of benefitting computational methods in order to get to know different life forms and their interactions. Thanks to recent technological developments in bioinformatics, analyzing and interpreting the biological data is now faster and more accurate. Identifying the different characteristics and structures of the human proteins and pathogen proteins that infect us, become essential to be more proactive in humans fight with pathogens. Protein sequence analyzing makes it possible to predict various potential interactions between virus and human proteins by identifying the significant features and classifying them with respect to their similarities and differences.

In this project, Virex, the objective is to develop a model in order to predict possible infections between human and virus proteins through gathering protein sequences by web scraping, determining the main features of those proteins and classifying the various interactions by using machine learning algorithms.

The developed model consists of three subsystems: Data Gathering, Feature Engineering, and Prediction. While Data Gathering subsystem is responsible from extracting the protein sequence data from the web, Feature Engineering subsystem converts that data into required form and identifies the significant features in predicting the interactions between human and virus proteins which is accomplished by Prediction subsystem.

The model is tested with different data sizes and using different prediction algorithms. The results obtained with the proposed model, Virex, are really promising as 0.97 accuracy rate is obtained, outperforming the accuracy values provided in the literature.

Implementation Details

Virex is implemented in **Python language** on **Google Colab environment**.

Initial Pathogen-Host Interaction data was downloaded from **Phisto DB web page** then imported into the project with the help of **Pandas** library. Data visualizations are made with **Matplotlib** module.

Custom URLs are created with the path of unique virus and human proteins ids. With the help of **Urllib** library, the sequence data on **Uniprot.org** is extracted. In order to decrease the execution time in web scraping, **multiprocessing** library is used and 5 different processes work at the same time. Hardware –software mapping is visualized as it can be seen on Figure A. The resulting performance upgrade can be seen on Table 1 which was calculated with **time** library. Then human and virus sequences are **cleaned** from the unnecessary characters and saved into a CSV file. Some of the sequence data which are not available on the web is written on a txt file with **pickle** library to be excluded on all records later on. Track of the infection between the virus and human proteins are kept with the IDs on the initial data so the gathered sequences are joined according to those IDs. **Cartesian product** between the unique human and virus proteins are created with **Itertools** library just for creating the infection negative records. By comparing with the positive infection records negative infections are labeled as Infection False. Infection column is stated as the target column, Virus Sequence and Human Sequence columns are stated as the features.

Amino acid sequences are split into their amino acids and counts of the characters are listed as a matrix by using **Count Vectorizer**. Then by using **Tf-idf Transformer**, their term frequencies and inverse document frequencies are calculated and their product is written on the related columns. The frequencies of matching virus and human amino acids are added together and divided into 2 in order to have their vectoral frequency representation as a match. Infection column is binarized by using **Label Binarizer**.

Features and labels are split into training and test sets with **Sklearn** library (0.2 rate) For experimental purposes, Virex applies **different portions of virus-host match data** into several ML algorithms to predict the infection between the virus and host proteins.

Those algorithms are as follows:

- Random Forest Classifier
- Naïve Bayes Classifier
- Support Vector Classifier

The classification algorithms, confusion matrices, classification reports are implemented with Sklearn library.

Execution Time	Technique
0.20561695098876953 seconds	Saving and reading from files
788.8343715667725 seconds	Web scraping with urllib (Multiprocessing)
4148 seconds	Web scraping with urllib

Table 1. Execution times according to the data collecting method

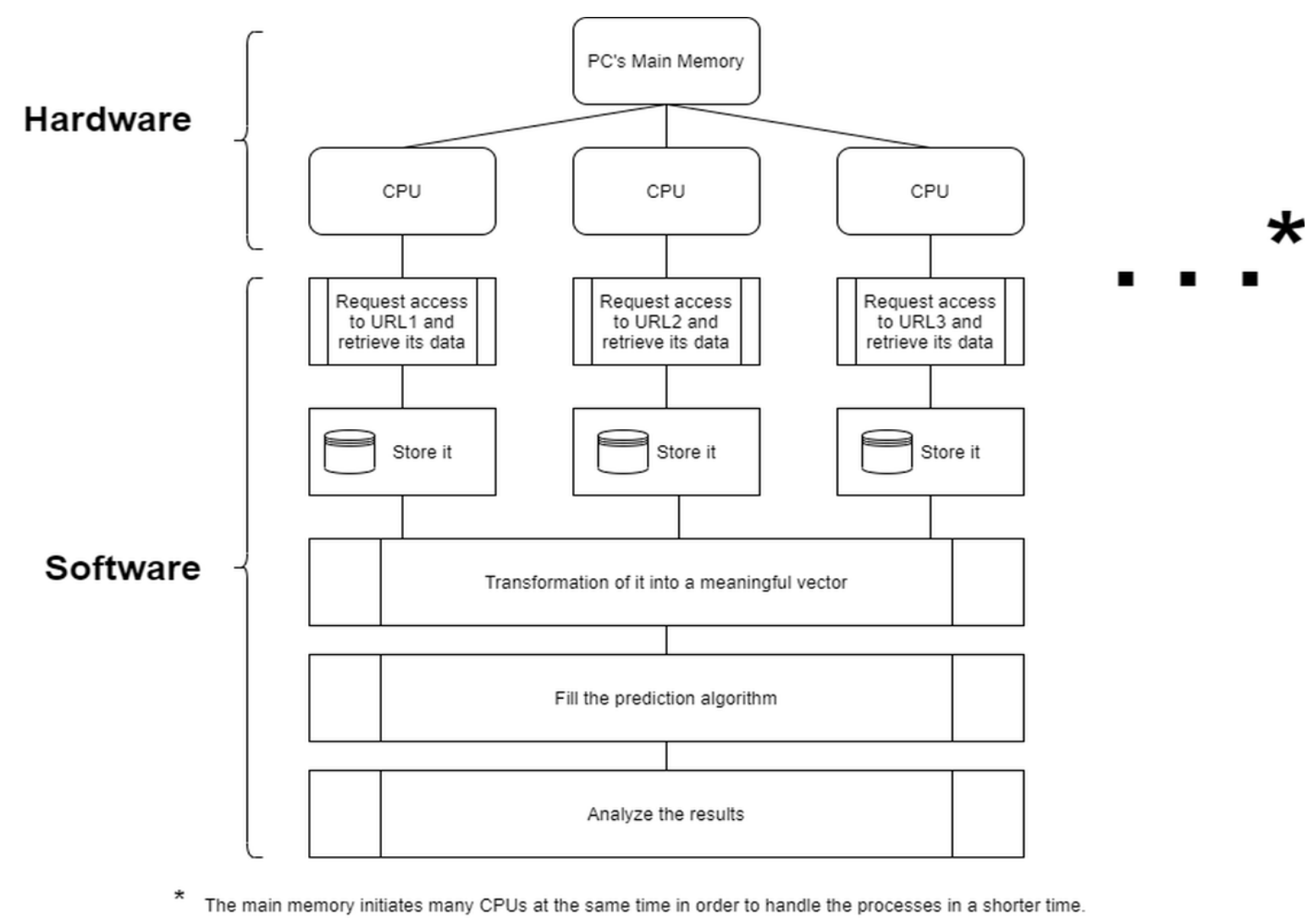


Figure A. HW-SW Relations

Results

HumanProtein	VirusProtein	
P15311	Q96P98	
P28799	Q97072	
P62937	P88761	
P42937	P88585	
P23284	Q77161	
1123888	K7PFA8	Q88851
1123889	K7PFA8	P80036
1123890	K7PFA8	Q88851
1123891	K7PFA8	Q88851
1123892	K7PFA8	Q70663

Figure 1. 11 Million records of virus-host sequence matches

	precision	recall	f1-score	support
0	0.85	0.93	0.88	11975
1	0.87	0.75	0.81	8025
accuracy			0.86	20000
macro avg	0.86	0.84	0.85	20000
weighted avg	0.86	0.86	0.85	20000
0.85555				

Figure 3. Prediction of 100,000 rows of data on Support Vector Classifier

	precision	recall	f1-score	support
0	0.97	0.86	0.91	2080
1	0.97	0.99	0.98	7920
accuracy			0.97	10000
macro avg	0.97	0.93	0.95	10000
weighted avg	0.97	0.97	0.97	10000
0.966				

Figure 5. Prediction of 50,000 rows of data on RF Classifier

	precision	recall	f1-score	support
0	0.97	1.00	0.99	192143
1	0.92	0.33	0.48	7857
accuracy			0.97	200000
macro avg	0.95	0.66	0.74	200000
weighted avg	0.97	0.97	0.97	200000
0.972545				

Figure 7. Prediction of 100,000 rows of data on RF Classifier

	a	c	d	e	f	g	h
0	0.275332	0.048583	0.229607	0.429746	0.127160	0.224107	0.121395
1	0.280314	0.414100	0.282211	0.170795	0.077886	0.254401	0.141512
2	0.389209	0.047461	0.136771	0.307771	0.134229	0.377830	0.121963
3	0.237864	0.082497	0.138874	0.322872	0.289529	0.365640	0.064423
4	0.271197	0.051562	0.186108	0.267386	0.189068	0.378100	0.060266
499995	0.289511	0.068415	0.139566	0.281287	0.139598	0.377536	0.104068
499996	0.235568	0.130765	0.283510	0.091736	0.372346	0.050284	0.050284
499997	0.307043	0.075457	0.203418	0.207805	0.157793	0.307131	0.090176
499998	0.284705	0.101579	0.164418	0.139286	0.188887	0.482547	0.099710
499999	0.224623	0.089759	0.223220	0.116038	0.216241	0.389128	0.070611

Figure 2. Tf-idf value averages of Human and Virus Sequence Vectors

	precision	recall	f1-score	support
0	1.00	1.00	0.75	11975
1	1.00	0.00	0.00	8025
accuracy			0.60	20000
macro avg	0.80	0.50	0.37	20000
weighted avg	0.76	0.60	0.45	20000
0.5989				

Figure 4. Prediction of 100,000 rows of data on Naive Bayes Classifier

	precision	recall	f1-score	support
0	0.93	0.98	0.95	11975
1	0.96	0.89	0.92	8025
accuracy			0.94	20000
macro avg	0.94	0.93	0.94	20000
weighted avg	0.94	0.94	0.94	20000
0.93955				

Figure 6. Prediction result of 100,000 rows of data on RF Classifier

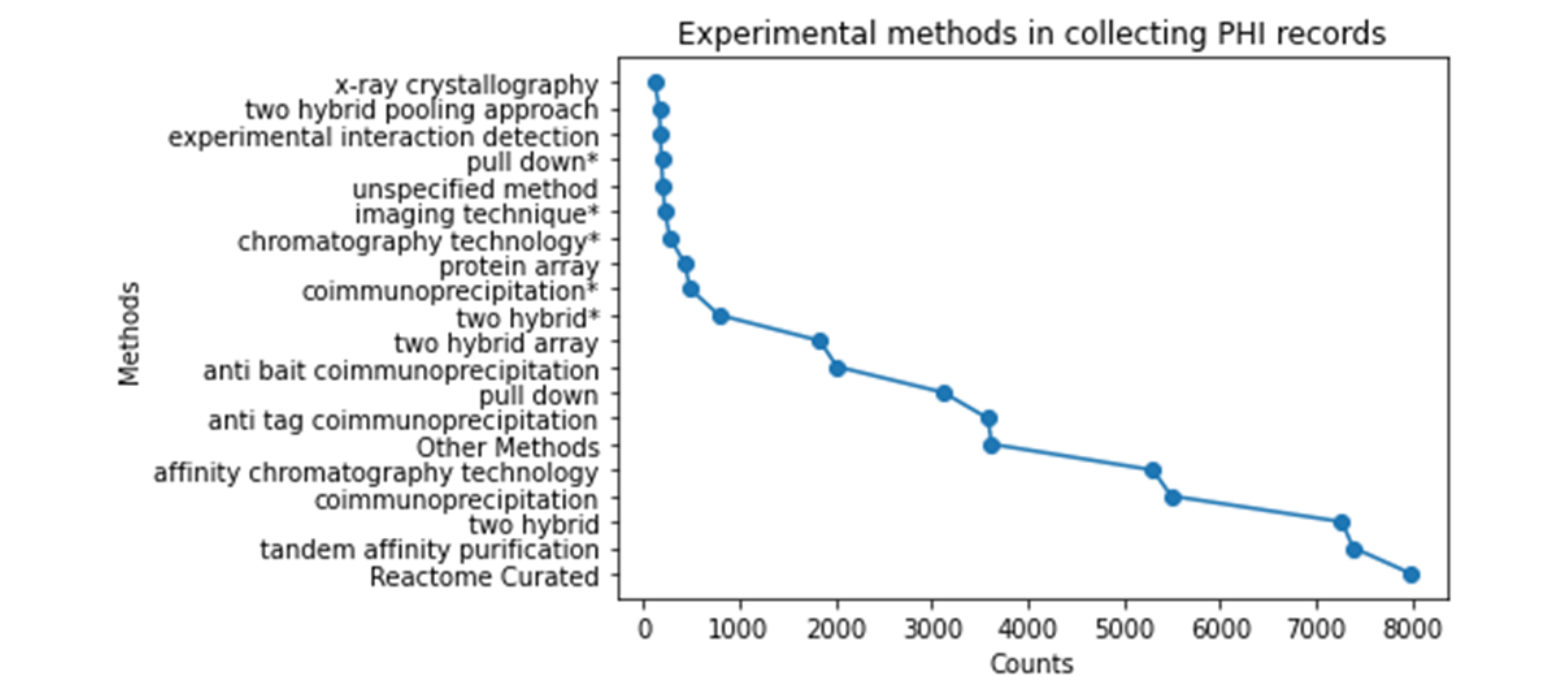


Figure 8. Experimental methods used for interaction observation in PHISTO data

Conclusion

This project aims at predicting human and virus protein interactions proposing machine learning based models. The developed system composed of three subsystems being Data Gathering, Feature Engineering, and Prediction is implemented with various data sets and the results show that Virex is capable of successfully finding and extracting data from the internet, processing the data through feature engineering and predicting the possible infections.

-RF Classifier: Best results are performed with RF algorithm among all the other performed algorithms. If we evaluate different data portions, the accuracy of the classification was the highest with 1M row of data (0.97). 50K rows of data resulted as the next best accuracy (0.96). Then as the third experiment 100K rows of data resulted with an accuracy of 0.939.

For the 1M rows: 191929 of the infection negative values are guessed correctly, only 214 of them are wrong labeled as infection positive values. 2580 of the infection positive values are labeled correctly, 5277 of the positive infection values are mislabeled as negative infection.

For the 100K rows: 7114 of the infection positive values and 11677 of the infection negative values are classified correctly. 298 of the infection negative values are misclassified as infection positive, 911 of the infection positive values are misclassified as infection negative.

For the 50K rows: 7862 of the infection positive values, 1798 of the infection negative values are correctly labeled. 282 of the infection negative values are misclassified as infection positive, 58 of the infection positive values are misclassified as infection negative.

-SVM Classifier: Even though SVM was slightly underperforming than the RF classifier, despite that it gave a high accuracy which was 0.855. It guessed 11105 infection negative values and 6006 infection positive values correctly. Only it misunderstood the 870 of the infection negatives and 2019 of the infection positives.

-Naive Bayes Classifier: NB algorithm was the worst in accuracy among all the other algorithms. From the confusion matrix and classification report it can be seen that the algorithm labeled 8022 of the infection positive values as infection negative values. That's why this accuracy is smaller than the others. But it guessed all of the infection negative values (11975), and 3 of the infection positive values correctly so the accuracy slightly increases. There are 0 mislabeled infection negative values.

Bibliography

- Goede Gökyokuş, Thesis Report Virex, Işık University, Jan 2021
- Sklearn, Count Vectorizer, lastly yok. <https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html>
- Label Binarizer. <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.LabelBinarizer.html>
- MultinomialNaiveBayes. <https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html>
- Train test split. <https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html>
- Tf-idf Transformer. <https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfTransformer.html>
- Google, Collaboratory. <https://research.google.com/collaboratory/faq.html>
- Pandas, lastly yok. <https://pandas.pydata.org/>
- Python Internet Access using Urllib, Request and Urlopen. <https://www.guru99.com/accessing-internet-data-with-python.html>
- Python, Itertools. <https://docs.python.org/3/library/itertools.html>
- Multiprocessing. <https://docs.python.org/3/library/multiprocessing.html>
- Pickle. <https://docs.python.org/3/library/pickle.html>
- Time. <https://docs.python.org/3/library/time.html>
- Yin, Tony, Random Forest. 12/06/2019. <https://towardsdatascience.com/understanding-random-forest-58381e0602d2>
- Matplotlib. <https://matplotlib.org/>
- UnProt. 24/05/2020. <https://www.uniprot.org/w/index.php?title=UnProt&oldid=958490101>
- Support vector machine. 20/10/2020. <https://en.wikipedia.org/w/index.php?title=Support_vector_machine&oldid=984421646>
- Train test split. <https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html>
- Ilknur Karadeniz, Salih Durmuş Tekir, Tunahan Çakır, Emre Arda, Ali Semih Sayınbaş, Gökhan Konuk, Mihai Konuk, Haaret Sarıyer, Azat Uğurlu, Arzuhan Özgür, Fatih Erdoğan Sevilgen, Katlı Ö. Ülgen. #PHISTO: pathogen-host interaction search tool. Bioinformatics. 20(2013): 1357-1358. <https://doi.org/10.1093/bioinformatics/btt127>
- UnProt. 24/05/2020. <https://en.wikipedia.org/w/index.php?title=UnProt&oldid=958490101>
- Kohl, Shivam, Understanding a Classification Report For Your Machine Learning Model. 18/11/2019. <https://medium.com/@kohlishivam522/understanding-a-classification-report-for-your-machine-learning-model-88815e2ce397>