
REPORT ASSIGNMENT 2

CHEMINFORMATICS WINTER SEMESTER 2022-23

Stephan Liu
Universität Tübingen
stephan.liu@student.uni-tuebingen.de

November 7, 2022

ABSTRACT

KNIME, also called Konstanz Information Miner, is an open-source data analytics program with the goal to advance the impact and understanding of data science through visual programming.¹ Thus, this program makes abstract tasks easier to understand as the problem is visualized. The graphical user interface is simple and clear and makes operations fairly user friendly.

1 Introduction

The importance of data and data sets cannot be neglected in science. The goal is to process these data in such a way that noise is minimized and still representative enough to be useable. The sheer volume, variance, and dimensionality of the data available make it difficult for the individual to identify, extract and interpret them correctly, which increases the importance of data analysis and statistical learning significantly.² Data analytics and machine learning simplify and advance our understanding in our endeavor to work with data.

Natural to say as the amount of data increases it gets gradually harder to keep clarity over the data and harder to reprocess the data. Accordingly, there are data processing programs that are used to maintain better uniformity. Besides well-known software like Rapidminer, IBM SPSS statistics and SAP analytics cloud, there are free open-source data analytics like KNIME that provide easy access for university purposes. In the field of science KNIME can be used to process large amounts of raw data such as a high-throughput-screening (HTS) produces.

2 Material and Methods

For the first part of the task three chemical structures were provided and their free ion pairs, as well as the hybridizations of substructures a) and c), were displayed.

To analyze the HTS data a KNIME workflow was created.

The 500 CSV files produced by the HTS contains raw data that are all read in into one single table. All screening plates are then visualized using heat maps and checked for their quality using the Z-prime factor considering the positive and negative controls.

The plates then were divided into plates of good quality plates with a Z-score of 0.5 to 1.0 and plates of bad quality with a Z-score of 0 to 0.5. Hit picking using the normalized percent inhibition (NPI) normalization with a $NPI \geq 70\%$ was used for the good plates. Hit picking using percent of control (POC) with a $POC \geq 60\%$ was used for the bad plates.

The data of the picked hits then was merged again and joined with the screening library file. The matches between the compound data and screening library file then were identified and its properties filled into a new CSV file.

¹<https://en.wikipedia.org/wiki/KNIME>

²<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8511823/>

3 Results

Figure 1 shows the chemical structures with free ion pairs and their hybridizations.

As for structure a) the upper all C-atoms have a sp^2 -hybridization. Both N-atoms have one free ion pair and therefore a sp^2 -hybridizations. The lower ring has the same properties as the upper structure with the difference that the NH^+ has no free ion pairs but an sp^3 -hybridization.

Structure b) has one free ion pair at the central C-atom as it has only three bonds to its neighbors.

Structure c) shows the upper molecule with two O-atoms that have their two free ion pairs drawn in. They both have an sp^2 -hybridization.

Meanwhile the lower substructure contains an O-atom with a sp^2 -hybridization and two free ion pairs and another O-atom with a sp^3 -hybridization and three free ion pairs.

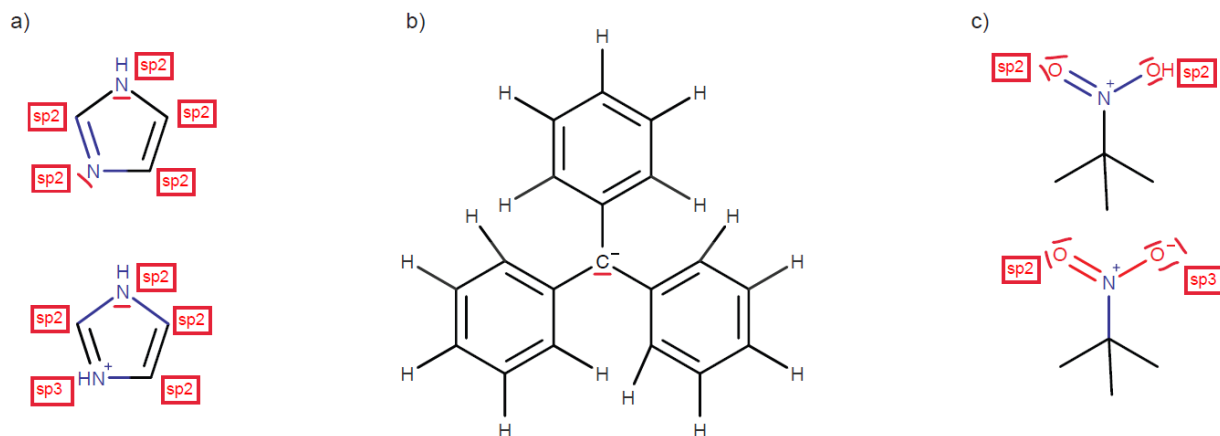


Figure 1: Chemical structures with free ion pairs and hybridizations.

Figure 1 shows the heat map that is generated using HCS tools plate heat map viewer therefore every plate shows what can be considered good or not within the plates itself in column one.

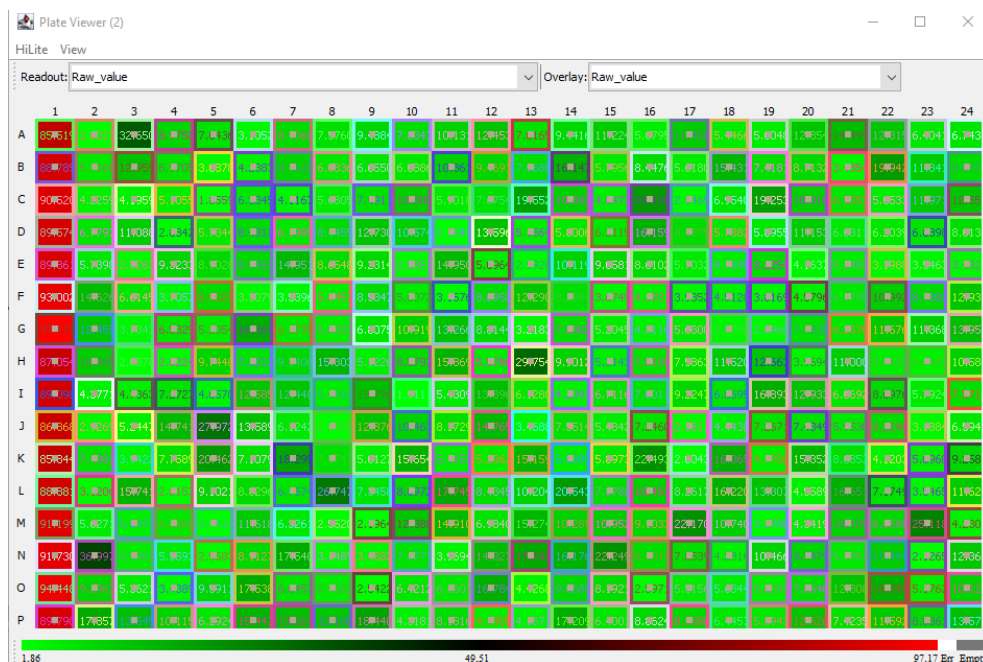


Figure 2: Heat map of plate 2 shown exemplary.

Figure 3 shows the output after merging the compound data with the screening library file and running through the node "Molecule to CDK" which allows to convert the SMILES properties to visualized molecular structures.

▲ Parsed molecules - 324 - Molecule to CDK (convert SMILES)

File Edit Hilite Navigation View

Table "default" - Rows: 194736 Spec - Columns: 12 Properties Flow Variables

Row ID	SMILES	Plate_ID	Plate_R...	Plate_Col	Sample	Project	Plate_J...	Plate_R...	Plate_C...	Sample...	Raw_y...	Raw_v...
Row0		A	A	B	TN001-01-02	HTS_Cel_Screen_FP	A	A	B	TN001-01-02	4.167	-0.779
Row1		A	A	C	TN001-01-03	HTS_Cel_Screen_FP	A	A	C	TN001-01-03	5.846	-0.417
Row2		A	A	D	TN001-01-04	HTS_Cel_Screen_FP	A	A	D	TN001-01-04	16.239	1.823
Row3		A	A	E	TN001-01-05	HTS_Cel_Screen_FP	A	A	E	TN001-01-05	9.636	0.4
Row4		A	A	F	TN001-01-06	HTS_Cel_Screen_FP	A	A	F	TN001-01-06	9.247	0.316
Row5		A	A	G	TN001-01-07	HTS_Cel_Screen_FP	A	A	G	TN001-01-07	11.768	0.859

Figure 3: Table with converted SMILES structures.

Figure 4 shows the output that is finally saved into the CSV file. It contains the header (1)sample ID, (2)molecular weight, (3)number of h-bond donors, (4) number of h-bond acceptors and (5)SlogP. It reports the chemical properties of the compounds of the screening library that are matching with the sample ID of the hits of the plates respectively.

chin-a2-Liu-Stephan.csv - Editor

Datei Bearbeiten Format Ansicht Hilfe

"Sample"	"Hydrogen Bond Acceptors"	"Hydrogen Bond Donors"	"Molecular Weight"	"XLogP"
"TN001-01-02"	1	1	162.06807956	2.823
"TN001-01-03"	3	2	176.094963004	0.9079999999999999
"TN001-01-04"	4	1	200.15247788	1.0899999999999999
"TN001-01-05"	3	2	129.078978592	-2.338
"TN001-01-06"	2	1	147.079647288	0.8039999999999998
"TN001-01-07"	4	3	131.058243148	-3.46
"TN001-01-08"	2	1	111.068413908	-0.2300000000000001
"TN001-01-09"	4	2	166.050905272	0.37199999999999994
"TN001-01-10"	4	1	178.004623776	0.8380000000000004
"TN001-01-11"	4	1	228.183778008	1.8060000000000003
"TN001-01-12"	4	3	126.05416081199999	-1.043
"TN001-01-13"	3	3	221.04865277599998	1.778
"TN001-01-14"	4	1	214.168127944	1.285
"TN001-01-15"	4	1	200.15247788	1.09
"TN001-01-16"	1	2	151.099714036	1.254
"TN001-01-17"	3	1	95.04834715999999	-0.5710000000000001
"TN001-01-18"	2	1	144.068748256	1.308
"TN001-01-19"	4	3	119.058243148	-1.698
"TN001-01-20"	3	2	165.078978592	0.5710000000000002
"TN001-01-21"	4	1	228.183778008	1.643
"TN001-01-22"	2	2	201.996761524	1.2200000000000002
"TN001-01-23"	4	3	144.089877624	-3.4400000000000004
"TN001-02-02"	2	1	122.08439831999999	0.20700000000000007
"TN001-02-03"	2	1	158.08439832	1.247
"TN001-02-04"	3	2	123.07964728799999	-1.105
"TN001-02-05"	4	2	214.168127944	1.7009999999999998
"TN001-02-06"	3	2	141.042593084	-3.078
"TN001-02-07"	4	2	188.035255208	-2.515

Zeile 1, Spalte 1 100% Windows (CRLF) UTF-8

Figure 4: CSV file with the chemical properties of the matching compounds.

References

- [1] Berthold, Michael R. and Cebron, Nicolas and Dill, Fabian and Gabriel, Thomas R. and Kötter, Tobias and Meinl, Thorsten and Ohl, Peter and Thiel, Kilian and Wiswedel, Bernd. KNIME - the Konstanz Information Miner: Version 2.0 and Beyond *SIGKDD Explor. Newsl.*, 2009.