

Manual of
SPCI (structural and physico-chemical interpretation) open-source software
version 0.1.5

Version (date)	Changes and comments
0.1.0 (02.02.2015)	<p>Changes from alpha version:</p> <ol style="list-style-type: none"> 1. More precise default SMARTS was added. 2. Cross-validation calculation was sped up and intermediate predictions are saved in text files. 3. Compounds, which cause errors in calculation of atomic properties with Chemaxon cxcalc tool, are excluded from further modeling. 4. Intermediate results of fragments contributions calculation are saved to a text file.
0.1.1 (07.02.2015)	<ol style="list-style-type: none"> 1. Fixed errors in text output of intermediate results of fragments contributions calculation. 2. Fixed error in loading of file with descriptors on 32-bit platforms.
0.1.2 (13.05.2015)	Fixed error in loading of a file with descriptors on 32-bit platforms.
0.1.3 (22.07.2015)	Added two automatic fragmentation schemes: detection of i) all rings and ii) Murcko frameworks.
0.1.4 (21.01.2016)	<ol style="list-style-type: none"> 1. Added automatic fragmentation scheme based on SMARTS to cleave bonds. 2. Changed routine of calculation of simplexes with hydrogen bonding labels. Simplexes with no labels A (acceptor), D (donor) or AD (acceptor and donor) are discarded. 3. Added support of multiple modeling properties in a single sdf-file (if property is missing put NA and such compounds will be discarded during modeling). 4. Default format of descriptors files is svm now (more effective to store sparse data). 5. GUI changes: <ul style="list-style-type: none"> - 10 last paths of loaded sdf-files are stored and can be easily loaded again - user may specify the number of cores to use for model building (default value is max_cores - 1) <p>Note. Backward compatibility is broken, old projects cannot be opened with this new version due to changes in project structure.</p>
0.1.5 (09.05.2016)	<ol style="list-style-type: none"> 1. Predict module (the tab "Predict") was added to predict properties of external datasets with the built models. 2. Project structure was changed: intermediate files required for calculation of non-user contributions are stored in the main project dir. 3. SiRMS submodule was updated. 4. Verbosity of output messages was decreased. 5. Many bug fixes (e.g. fixed calculation of specificity in classification models which crash the program).

Content

Overview	3
Concept and citations.....	3
Structural (<i>doesn't require Chemaxon</i>) and physico-chemical interpretation (<i>Chemaxon required</i>)	3
Installation and launch.....	3
Windows.....	4
Linux (Ubuntu):.....	5
Workflow:	7
Step 0. Data preparation and project start.	7
Step 1. Build models.....	7
Step 2. Calculation of the fragments' contributions.	8
Step 3. Plot contributions and save plot in png file.	8
Alternative step 3.	9
Step 4. Predict properties of new compounds with built QSAR models	10

Overview

The SPCI software was designed for (semi)automatic extraction of structural features and their contributions to an investigated property from chemical datasets. It's a bunch of Python3 scripts. In order to simplify the usage a GUI was developed. It has a limited number of options but suits well for most needs.

Concept and citations

The idea of structural and physico-chemical interpretation along with comparison with MMP approach is disclosed in the following references. Please cite them if you use this software.

- 1) Polishchuk, P. G.; Kuz'min, V. E.; Artemenko, A. G.; Muratov, E. N. Universal Approach for Structural Interpretation of QSAR/QSPR Models. *Molecular Informatics* **2013**, 32, 843-853.
- 2) Polishchuk, Tinkov, Khristova, Ognichenko, Kosinskaya, Varnek, Kuz'min Structural and physico-chemical interpretation of QSAR models, **2016** (in preparation).

Structural (*doesn't require Chemaxon*) and physico-chemical interpretation (*Chemaxon required*)

Structural interpretation returns only overall fragments contributions while physico-chemical interpretation can additionally estimate contributions of some physico-chemical factors (electrostatic, hydrophobic, hydrogen bonding and dispersive terms). Installed Chemaxon is required for physico-chemical interpretation.

If you have a license for Chemaxon check whether it is installed and PATH variable is correctly configured (to launch `standardize` and `cxcalc` from command line). If `standardize` or `cxcalc` command are not recognized from your command line, add JChem bin folder to PATH variable.

Installation and launch.

There are some dependencies which should be installed:

	Tested environments	
	Windows	Linux
Python	>= 3.2	3.4
matplotlib	1.4.3	1.3.1
numpy	1.9.2	1.10.4
scipy	0.15.1	0.16.1
scikit-learn	0.16.1	0.17
indigo toolkit	1.1.12	1.2.1

Windows

You may install all Python dependencies by yourself and download the latest version of the application from github in order to install and upgrade it (see Linux installation section for some details).

To make life easier for Windows users WinPython distributions (32 and 64-bit) were prepared and zipped. They contain all required dependencies for immediate start.

Step 1. Visit the page http://qsar4u.com/pages/sirms_qsar.php, download a desired WinPython distribution and unzip the archive into a folder.

Step 2. Download spci.zip and unpack into the same folder where WinPython was unpacked. You should obtain an installation folder with three items: spci and WinPython folders and start-spci.bat.

Program_installation_folder

```
| -spci/  
| -WinPython/  
| -start-spci.bat
```

To launch the application choose start-spci.bat.

Update:

To update the application just remove spci folder from the installation folder and unzip newly downloaded spci.zip in it.

Linux (Ubuntu):

Install `numpy` and `scipy` packages from Ubuntu repositories or with `pip3`:

```
sudo apt-get install python3-numpy python3-scipy
```

If these packages were installed earlier I would recommend to upgrade them:

```
sudo pip3 install -U numpy
```

```
sudo pip3 install -U scipy
```

Run the following command to install missing libraries:

for Ubuntu 14.04:

```
sudo apt-get install libatlas-dev libatlas3gf-base
```

for Ubuntu 16.04:

```
sudo apt-get install libatlas-dev libatlas3-base
```

Then install `scikit-learn`:

```
sudo pip3 install scikit-learn
```

And make sure that atlas is used to provide the implementation of the BLAS and LAPACK linear algebra routines:

```
sudo update-alternatives --set libblas.so.3 \  
    /usr/lib/atlas-base/atlas/libblas.so.3
```

```
sudo update-alternatives --set liblapack.so.3 \  
    /usr/lib/atlas-base/atlas/liblapack.so.3
```

Download and install Indigo toolkit -

<http://lifescience.opensource.epam.com/download/indigo/index.html>

Simply copy all files and a folder containing in the supplied zip file to your `dist-packages` dir (or `site-packages` on Windows).

Then clone the repo:

```
git clone https://github.com/DrrDom/spci
```

Go into the created `spci` dir and init and update submodules:

```
git submodule init
```

```
git submodule update
```

To start the application call `python3 spci.py` from `spci` dir.

Update:

Update spci repository

```
git pull origin master
```

Update submodules if necessary

```
git submodule update
```

Workflow:

Step 0. Data preparation and project start.

For analysis a single sdf file is required, which contains compounds' structures and corresponding property value as a field in sdf file. Input sdf file may contain several modeling properties. The structures should be checked and standardized before run analysis (this is absolutely required in the case of only structural interpretation model, otherwise results can be distorted because standardization influence calculated descriptors).

As an alternative, you may prepare separate text file with compounds' property values (file must have a header and two columns compound name and property value separated by tab).

Place the sdf file in a separate folder which will be a project folder.

Step 1. Build models.

The screenshot shows the SIRMS QSAR interpretation software interface. The window has tabs for 'Build models', 'Calc contributions', and 'Plot contributions'. The 'Build models' tab is active. It contains several sections:

- 1. Structural/physico-chemical interpretation:** Radio buttons for 'Structural & functional (Chemaxon required)' and 'Structural only (no Chemaxon usage)'.
- 2. SDF with compounds:** A 'Path to SDF-file' dropdown and a 'property field name' dropdown.
- 3. Optional. Compound names. External text file with compound property values:** Radio buttons for 'Automatically generate compound names', 'Use compound titles from SDF file', and 'Use field values as compound names from SDF file'.
- 4. Models:** Checkboxes for various models: Regression (RF, GBM, SVM, PLS), Binary classification (0-1) (RF, GBM, SVM), Random Forest (RF), Support vector regression (SVR), Support vector classification (SVC), Gradient boosting regression (GBR), Gradient boosting classification (GBC), Partial least squares (PLS), and k-Nearest neighbors (kNN).
- 5. Number of cores to use:** A spinner box set to 3.

At the bottom are 'Build models' and 'Show statistics' buttons. A copyright notice '(c) Pavel Polishchuk 2014-2015' is at the bottom left.

*Red fields are required, green ones are optional.

1. Choose the type of analysis to perform.
2. Specify path to sdf file with compounds and choose the name of the field containing property values. Paths to 10 last opened sdf files are stored.
3. Optional. It is required if you want to save original compounds names or if property values are supplied in a separate file.
4. Choose the desired type of models. All listed models will be build with optimal parameters selected by grid search in 5-fold cross validation procedure.
5. Optional. The number of cores to be used during models building. By default all cores minus 1 will be used.

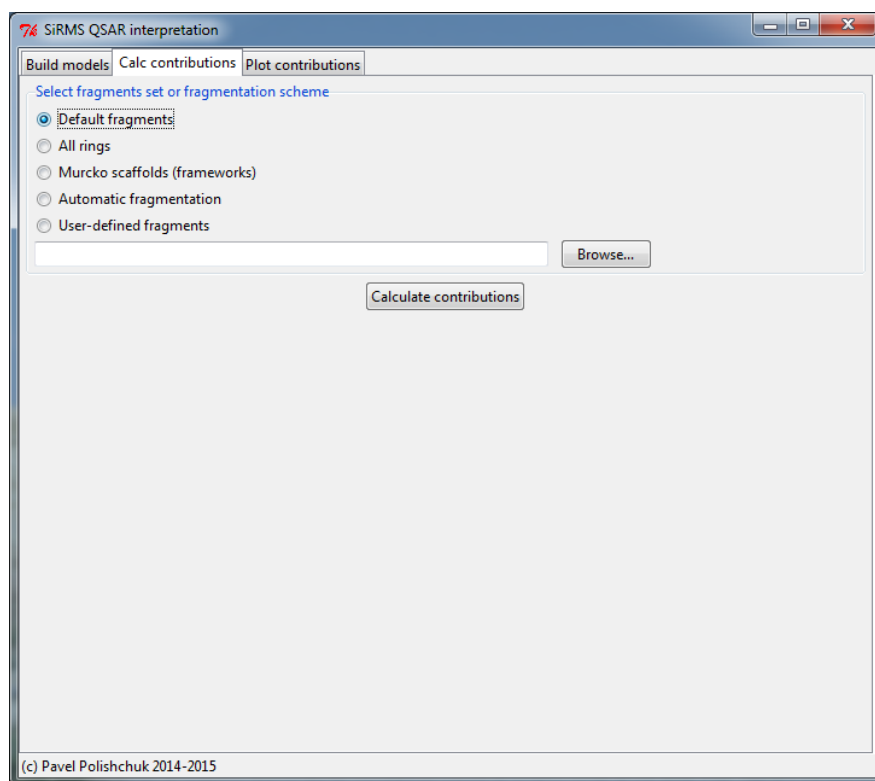
IMPORTANT:

To open an earlier created project you should choose the sdf file and property field name.

Step 2. Calculation of the fragments' contributions.

Contributions will be calculated only for those models specified on the previous tab.

If models had reliable predictive ability estimated by 5-fold cross-validation (Show statistics button) one may calculate fragments contributions.



default fragments – common functional groups and rings (specified in default.smarts file in program folder);

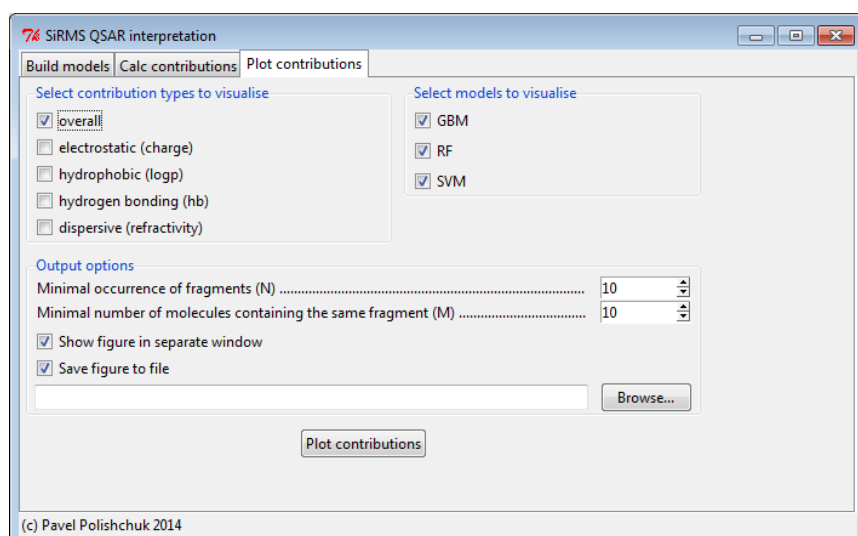
all rings – automatically detect all single and fused rings;

Murcko frameworks – automatically detect Murcko frameworks;

automatic fragmentation – split molecules on maximum three parts on bonds match SMARTS: [#6+0;!\$(*,#[!#6]))!@!=#[*]

user-defined fragments represented in SMILES or SMARTS notation, look at default.smarts as a file format example (Note: SDF usage was not tested).

Step 3. Plot contributions and save plot in png file.



After project was opened all available models will be listed in order to choose them for visualization.

Overall contribution (structural interpretation).

Other types is used for physico-chemical interpretation (to show the contribution of separate physico-chemical factors).

Alternative step 3.

You may create alternative visualization in order to customize plot output:

- 1) by yourself (parsing *_frag_contribution.txt files with calculated contributions located inside the modeling property folders and applying any of available tools);
- 2) by using the developed web-base tool which is suitable for visualization of relatively small set of fragments (number of columns in *_frag_contribution.txt should be less than 100k) otherwise it will be painful;

Full version - <http://158.194.101.252:3838/spci-vis/>

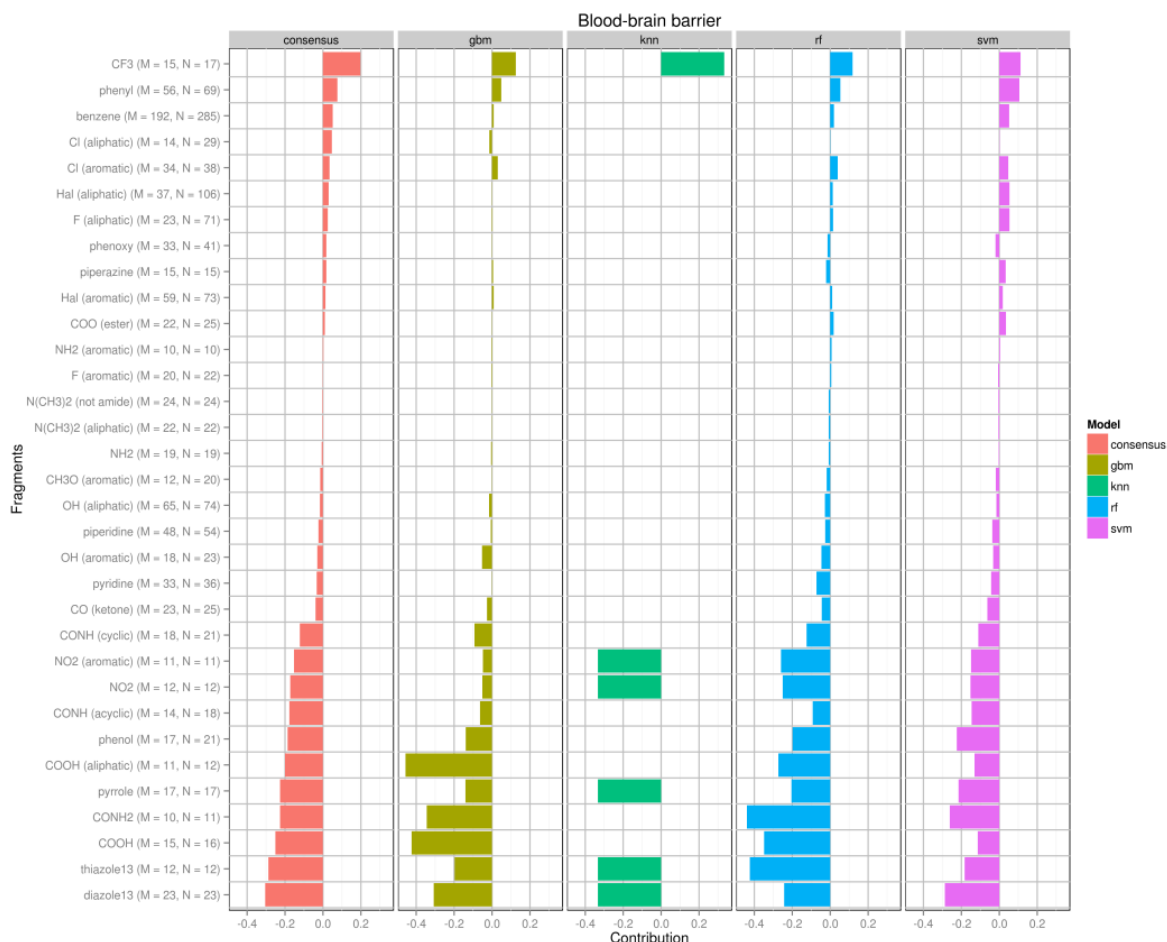
Demo version to just play with pre-loaded data - <http://158.194.101.252:3838/spci-vis-demo/>

- 3) by using `rspci` R package on your local machine. This package contains functions to facilitate data reading, modification, filtering and plot. It has more options than the web-based tool and recommended for advanced usage by R users.

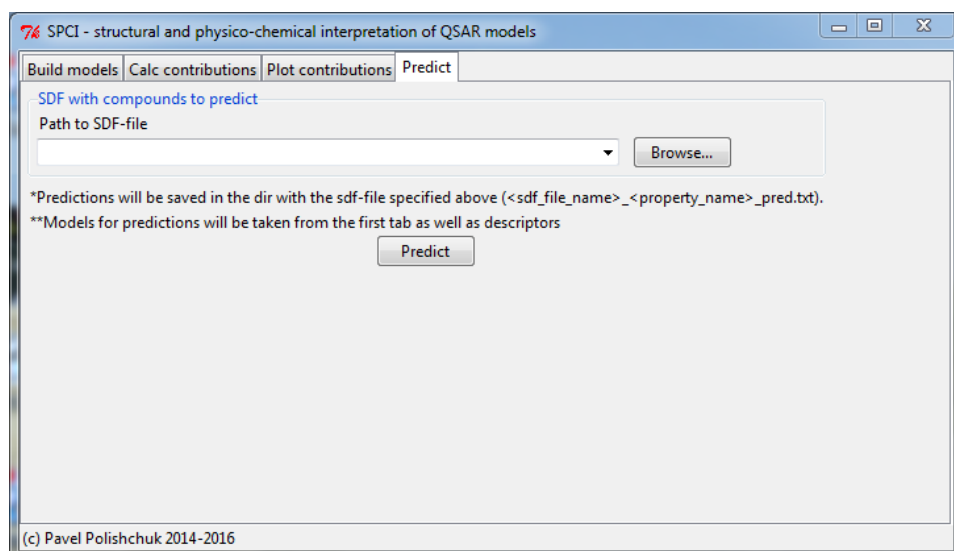
To install `rspci` package call from R console:

```
devtools::install_github("DrrDom/rspci")
```

Below is an example of contributions plot:



Step 4. Predict properties of new compounds with built QSAR models



Specify path to sdf-file with an external set of compounds and choose models and property for prediction from the first tab.

Applicability domain is estimated for each compound based on bounding box – if at least one descriptor has a value outside the range of this descriptor values for training set the compound is out of domain applicability.