**COmparative Molecular Field Analysis (COMFA)**

**By Gijs Schaftenaar**

### Overview of this tutorial

### Introduction

COMFA (COmparative Molecular Field Analysis) is a special case of QSAR.

QSAR (Quantitative Structure Activity Relationships) is linear regression technique based on data from known active molecules. QSAR can be applied when the 3D structure of the receptor is unknown. To apply QSAR, all that is needed are the activities, the 2D and/or 3D structures and properties/descriptors of the molecules. Of course, activities have to be measured, but 3D structures can be determined either by measurement (crystal X-ray analysis) or by calculation from the 2D diagram and (optionally) subsequent optimization.

COMFA is a QSAR approach where the descriptors consist of the Molecular Fields at each point in a 3D grid around the aligned set of 3D ligand molecules.

The aim of QSAR/COMFA is to derive a correlation between the biological activity of a set of molecules and their properties/Molecular Fields. Where some properties used with QSAR can be calculated when only the 2D structures of the molecules are available, other properties require a 3D structure. The Molecular Field used as descriptors with COMFA requires a set of 3D ligand structures, complete with hydrogens attached, aligned at their 'biologically active' conformations. Herein lies the problem with the COMFA method: since it is a ligand based method, it is preferably used at the beginning of a drug design project when no crystal structures of the target are available. Hence the 'biologically active' conformations of the ligands can only be guessed at. However what we can safely assume is that these conformations lie not too far above the optimal conformation in energy (at most 5 Kcal/mol) and that the conformations of the active ligands which have similar conformations are more likely to be valid.

How does COMFA work?

- COMFA uses a partial least-squares (PLS) analysis to predict activity from linear combinations of properties/descriptors.

What is needed before doing COMFA?

- Molecules with activities spanning about three log units of KI or IC50 values are required.
- The basic COMFA assumption is that similar molecules have similar activities, hence COMFA works best with a series of closely related molecules. As a consequence, the further the molecule, who's activity you are trying to predict with COMFA, is from the training set molecules, the worse your predicted activity will be.

In this tutorial, you will create a COMFA model and study its application.

### Setup the working environment

From the Unix shell (command prompt):

- Change directory to Practicals/07_COMFA/ by typing
  cd Practicals/07_COMFA/
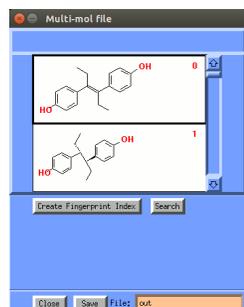- and call gmolden by typing
  gmolden

### The Molden Comfa interface

**Molden** provides an interface to the open-source package Open3DQSAR. Open3DQSAR uses molecular interaction fields (MIF) to calculate descriptors on each point on a 3D grid surrounding a set of pre-aligned molecules, supplied in the form of an **.sdf** file.

A series of 30 compounds with moderate to high activity for the estrogen receptor (ER-α) (*Endocrinology* **1997**, *138(9)*, 4022-4025) have been constructed and stored in the file **est+act.sdf**. This set, along with the measured activity data, will be used for training in a COMFA model.

First, we need to write out the activities stored in the **.sdf** file into a simple **.txt** file.
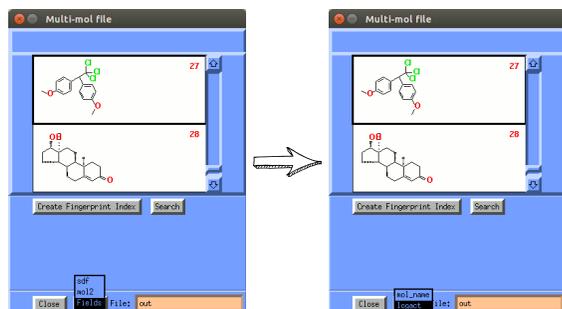
- In the molden control window, select **Read**. The *Molden File Select* window will pop up.
- Activate the *text entry field* below the **Filter:** string by clicking it. Now enter the filter **sdf** and hit the **Enter** key.
- Click the **est+act.sdf** file. The *Multi-mol file* window will pop up.

You should now have a *multi-mol file* with 30 molecules:



Take a look at the collection of molecules; You can copy a molecule from the *multi-mol file* to the molecular viewing area by single-clicking the 2D representation. Rotate the molecule by using the mouse (keep left mouse button pressed down). The majority of the molecules are steroids; do you recognize the steroid skeleton? Notice that there are also a number of non-steroidal molecules in the training set. Write down the numbers of the non-steroidal compounds. You will need them later on in the tutorial.
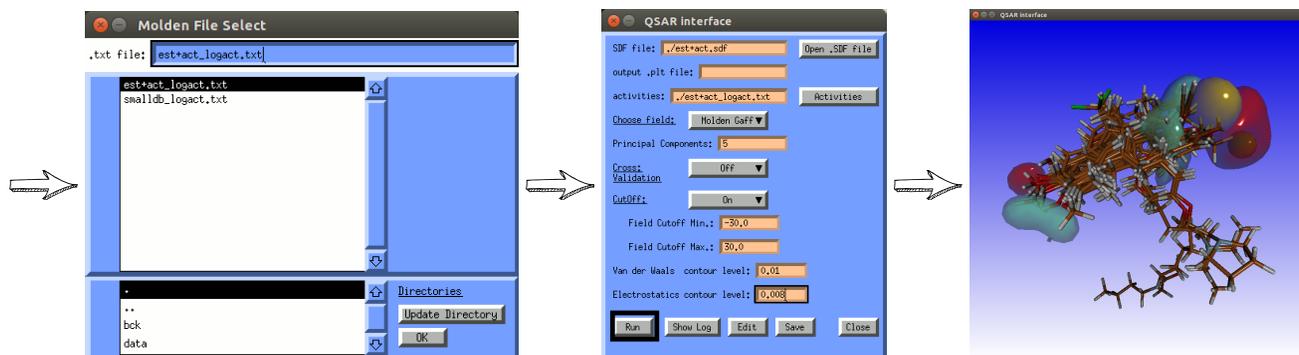
Now the activity data are exported from the *multi-mol file* to a text file, by clicking the **Save** button.:



The **est+act_logact.txt** text file has as first line *Affinity* and for each structure a line with the logarithm of the activity.

The activities are relative binding affinities, expressed as (nanomolar!) concentration values. In COMFA/QSAR the logarithms of concentrations are used, similar to using pH for the acidity of a solution. So the values we'll be using are -log(concentration*10⁻⁹).
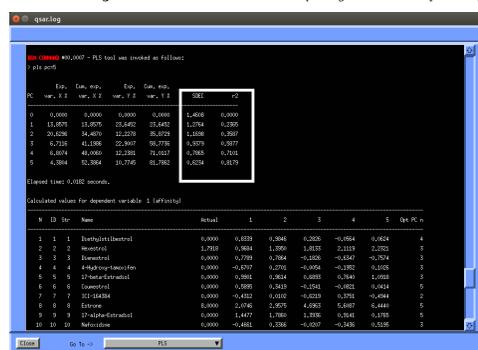
You should now have a **est+act.sdf** database with 30 molecules and a corresponding file **est+act_logact.txt** with activities.
Now let us start the molden *Open3DQSAR* interface:

For clarity's sake make the surfaces transparent by typing a **t** character in the display screen.

You will have to adjust the default values for the van der Waals and Electrostatics contour levels to **0.01** and **0.008** respectively. The higher these values are the more compacted the surfaces become. Higher contour values also mean greater correlation to the activity. The light blue and red surfaces show the correlation of the electrostatics part of the molecular field and the activity. These surfaces concentrate on the two OH groups which are essential for activity. Green contours indicate regions where a relatively bulkier substitution would increase the activity, whereas the yellow contours indicate areas where a bulkier substituent would decrease the activity. The regions where increasing the positive charges would increase the activity are in blue and regions where increasing the negative charges increases the activity are in red.

Hit the **Show Log** button to view the contents of the *qsar.log* file. This will open the *qsar.log* window. In this window select the **PLS** option to set the viewport to the PLS part of the file.



In the *Open3DQSAR* log file you will find **SDEC** (**S**tandard **D**eviation of **E**rror of **C**alculation) and $r^2$
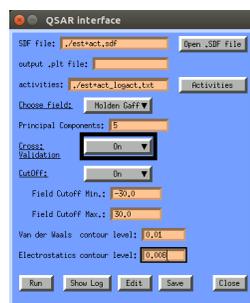
---

**Perform a cross-validated PLS analysis**

In a partial least-squares (PLS) analysis, two factors are important: the number of components used in the regression equation (which will be dealt with later) and the (usually squared) correlation coefficient. A *non*-cross-validated PLS analysis gives a squared correlation coefficient usually indicated by $r^2$. This number, which is also used in (multiple) linear regression, is between zero and one and expresses the quality of the PLS analysis. It indicates the proportion of the variation in the dependent variable (here the activity) that is explained by the regression equation and its value should be as close to one as possible. However, $r^2$ expresses the quality of the data fit rather than the quality of prediction (which is what we are actually interested in).
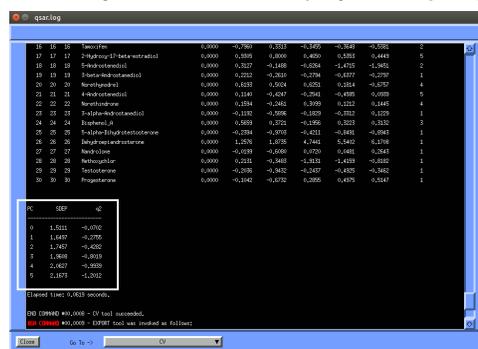
To express the predictive power of the analysis, the **cross-validated $r^2$**, usually indicated by $q^2$, is used. In cross-validation, one value is left out, a model is derived using the remaining data, and the model is used to predict the value originally left out. This procedure is repeated for all values, yielding $q^2$. $q^2$ is normally (much) lower than $r^2$ and values greater than 0.5 already indicate significant predictive power.

After the cross-validated PLS analysis, you will determine the optimal number of components. Recall that the components are linear combinations of the variables (of which you have very many!), ordered in such a way that the first component will describe most of the variation in the activity, the second most of the remaining variation, etc. The cross-validated PLS analysis will be carried out for different numbers of components. As a rule of thumb, the number of components should not exceed one-third of the number of molecules. More components would lead to a model that is *overtrained* - it has a better fit to the training data but the predictive power is diminished.

To perform a cross-validated PLS, we switch the *Cross Validation* button to **On**.



Hit the **Show Log** button to view the contents of the *qsar.log* file. This will open the *qsar.log* window. In this window select the **CV** option to set the viewport to the cross validation part of the file.



- In the *OPEN3DQSAR* log file you will find the Standard Error of Prediction (**SEP**) and $q^2$.

---

**Determine the optimal number of components**

The two important factors here are the **standard error of prediction** (SEP) and the **cross-validated $r^2$** (usually called $q^2$). SEP should be as low as possible and $q^2$ should be larger than 0.5 (lower values indicate poor predictive power).

- Find the number of components having the lowest SEP and note the corresponding $q^2$.
- To accomplish this, in the *QSAR nterface window*, specify the number of components under **Principal Components:**.
- Check whether an additional component increases $q^2$ with more than 5%.
- If so, the larger number should be used; otherwise, take the lower (original) number of components.

The resulting number of components will be used in the non-cross-validated analysis.

You will find that you can not get the $q^2$ above 0.5 with this combination of dataset and descriptors. In fact the $q^2$'s are negative !
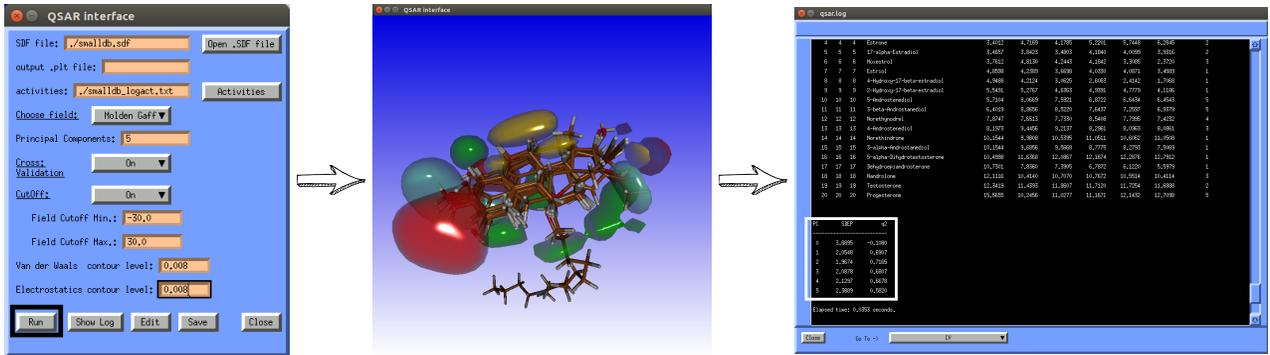We will try to get the $q^2$ above 0.5 by switching to a more focussed dataset: we will remove the non-steroidal compounds from the dataset.
Remember you wrote down the names of the non-steroidal compounds ? Now remove all non-steroidal compounds from the dataset:

- Go through the list of entries in the **est+act.sdf** file by hiting the *arrow down key* and hit the *Delete key* on any non-steroidal compound. When done deleting, click with the left mouse button, the *File:* text entry field and type **smalldb**. Next hit the **Save** button and select the *sdf* option.
  Next, click the **Save** button again, but select the *Fields* option. The **smalldb_logact.txt** file will have been created.

Now repeat the steps of COMFA model creation and cross-validation for the steroidal dataset and determine the optimal number of components. You will have to adjust the default values for the van der Waals and Electrostatics contour levels to **0.008** and **0.008** respectively.

**Determine the optimal descriptors and components**

By varying the number of components it is possible to get $q^2$ as high as 0.717.