

Worked Example:

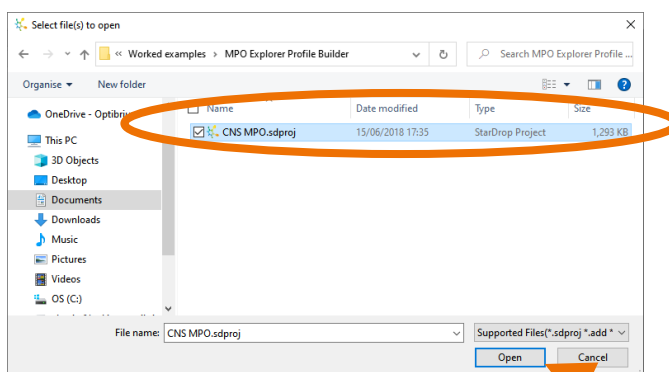
Automatically Building a Scoring Profile with Rule Induction

In this example, we will use the Profile Builder in StarDrop's MPO Explorer module to derive a multi-parameter scoring profile based on a data set initially described by Wager *et al.* [ACS Chem. Neurosci. 1 p. 435 (2010)]. The authors used this data set to develop a multi-parameter optimisation method for selection of compounds intended for CNS indications. The 'CNS MPO score' derived by Wager *et al.* is calculated as the sum of the values of desirability functions for six physicochemical parameters, calculated logP (clogP), calculated logD at pH 7.4 (clogD), molecular weight (MW), topological polar surface area (TPSA), number of hydrogen bond donors (HBD) and the pK_a of the most basic centre (pK_a), resulting in a value between 0 and 6. The authors compared the CNS MPO score for a set of 119 marketed drugs for CNS targets with 108 Pfizer CNS candidates and found that 74% of the marketed drugs achieved a desirability score of ≥ 4 compared with only 60% of the Pfizer candidates.

The scoring profile derived by MPO Explorer will contain one or more rules corresponding to combinations of properties that significantly increase the chances of identifying a drug, and we will compare these with the results of the CNS MPO score.

Exercise

- In StarDrop, open the project file **CNS MPO.sdproj** by selecting **Open** from the **File** menu.



The project contains two data sets, the first of which (called CNS MPO) contains 227 compounds, 119 drugs for CNS targets and 108 Pfizer development candidates that did not reach the market (published by Wager *et al.*).

For each compound, 6 properties have been previously calculated, labelled MW, CLOGP, TPSA, CLOGD, HBD and PKA, as described in the introduction. For comparison, the CNS MPO score, as defined in Wager *et al.*, is also included in the data set.

The screenshot shows the StarDrop - CNS MPO software interface. The main window displays a table with 11 rows of data. The columns are: Pub-Name, Set, PKA, HBD, CLOGD, and TPSA. The data is as follows:

| Pub-Name | Set | PKA | HBD | CLOGD | TPSA |
|----------------|------|------|-----|-------|------|
| 1 Acamprostate | Drug | 1 | 2 | -6.1 | 8 |
| 2 Alprazolam | Drug | 1.9 | 0 | 2.5 | 4 |
| 3 Amfebutamone | Drug | 7.2 | 1 | 3.3 | 2 |
| 4 Amisulpride | Drug | 9 | 3 | -0.2 | 10 |
| 5 Amphetamine | Drug | 9.9 | 2 | -0.6 | |
| 6 Aniracetam | Drug | 1 | 0 | 0.3 | 4 |
| 7 Apomorphine | Drug | 7.9 | 2 | 2.5 | 4 |
| 8 Aprepitant | Drug | 4 | 2 | 4.1 | 8 |
| 9 Aripiprazole | Drug | 6.7 | 1 | 5.5 | 4 |
| 10 Atomoxetine | Drug | 10.1 | 1 | 0.7 | 2 |
| 11 Bromazepam | Drug | 2 | 1 | 2.1 | 5 |

The interface includes a menu bar (File, Edit, View, Data Set, Tools, Custom Scripts, Help), a toolbar, and a sidebar with various tool options. The 'MPO Explorer' section at the bottom left has a 'Build profile...' button highlighted with an orange arrow.

We will use the data in the data set to find a scoring profile, based on these simple properties, with which to distinguish marketed CNS drugs from unsuccessful candidates.

- Click on the **Scoring** tab and, at the bottom, click the **Build Profile** button, under MPO Explorer, as shown above.

The **MPO Explorer** wizard, shown right, will be displayed enabling us to control the profile building process.

The screenshot shows the 'MPO Explorer' dialog box, specifically the 'Create Session' step. The 'Objective Type' is set to 'Category' (indicated by an orange arrow). The 'Set Split' is set to 'Automatic'. The 'Input Data' section shows 'Profile Name: Scoring profile - Set', 'Data Set: CNS MPO', 'Validation Set: <None>', and 'Test Set: <None>'. The 'Objective Column' is set to 'Set' (circled in orange). The 'Desired Outcome' is set to 'High'. The 'US Patent No. 9,367,812' is displayed at the bottom. The 'Next >' button is highlighted with an orange arrow.

- We would like to build a profile based on the categorical objective in the **Set** column. Therefore, choose **Category** under **Objective Type** and confirm that the **Objective Column** is defined as **Set**, as shown above.

Note that the **Desired Outcome** is **High**, i.e. we would like to maximise the value of the objective.

- Click the **Next** button to move onto the next page, **Define Categories**.

Here we can define the order of the categories, from the lowest at the top to the highest at the bottom. Remembering that we have chosen to maximise the objective, we need to order the categories such that the **Drug** is at the bottom.

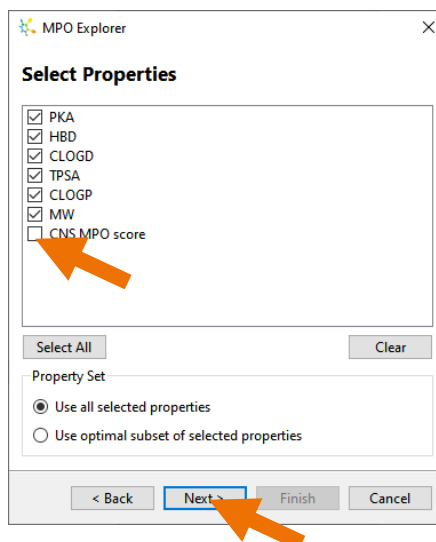
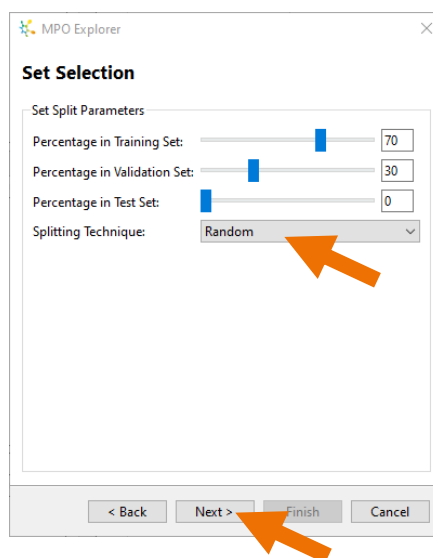
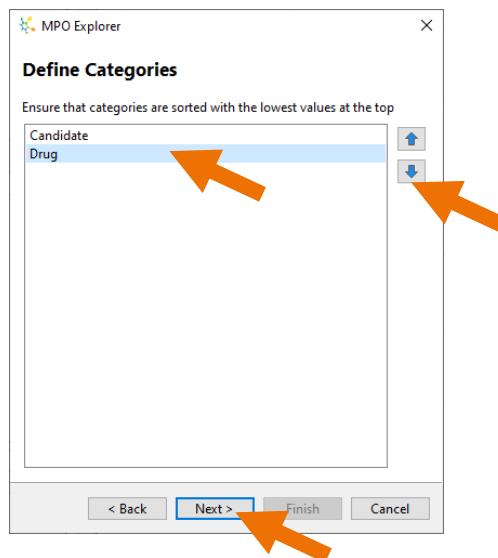
- Select the **Drug** category and click the down arrow.
- Click the **Next** button to move onto **Set Selection**.

Here we can define the parameters for splitting the data set into training, validation and test sets. In this case, we will use the defaults, putting 70% of the compounds into the training set and 30% into the validation set. Ideally, we would like an external test set, but given the small size of the data set, it is not practical in this case.

- Select **Random** as the **Splitting Technique**.
- Click the **Next** button.

On the **Select Properties** page (shown right) we can choose the properties that we would like to explore to identify a scoring profile. In this case, we would like to use the six simple compound properties, as described above.

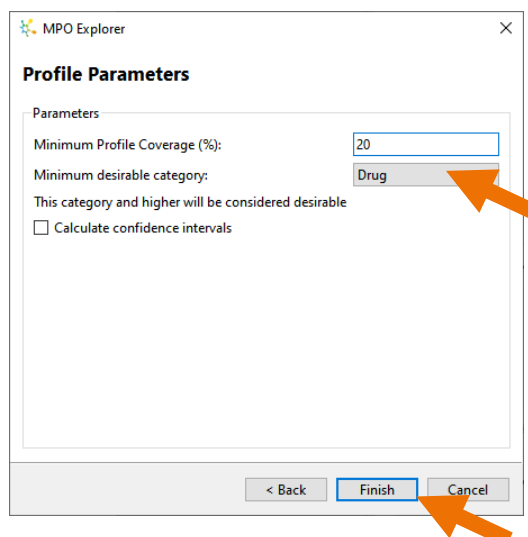
- Untick the **CNS MPO score** property, as shown right, to avoid using this.
- Click the **Next** button to move onto the last page of the wizard.



The **Profile Parameters** page enables us to specify some additional conditions. In this case, we will use the default minimum profile coverage of 20% (i.e. we will only consider rules that are applicable to $\geq 20\%$ of the compounds in the data set).

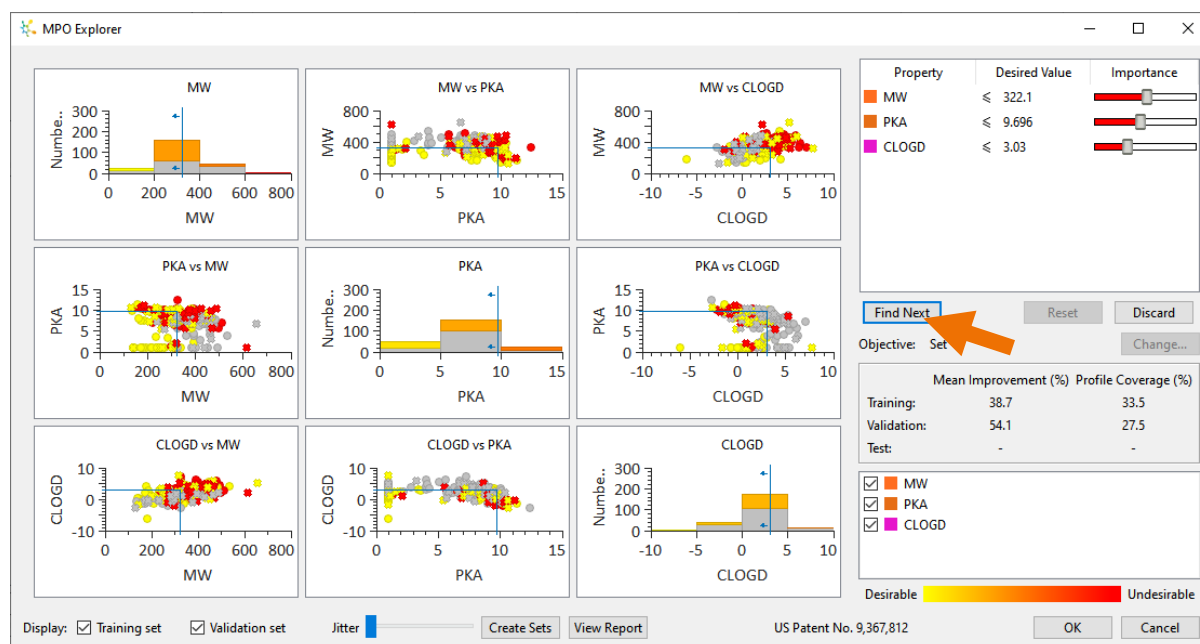
- Set the **Minimum desirable category** to **Drug**.

This is used to calculate the performance statistics for the rule's ability to distinguish desirable and undesirable outcomes. This parameter is most useful if the objective is a continuous property or has more than two categories.



- Finally, click the **Finish** button to begin the profile building session.

Once the first rule has been found, this will be displayed within the MPO Explorer profile analysis window, as shown below.



The rule is shown in the top-right. In this case, the rule suggests that compounds with a MW < 322.1 and a most basic PKA < 9.696 and a CLOGD < 3.03 will have an increased chance of success. The criteria for MW and PKA are slightly more important than CLOGD.

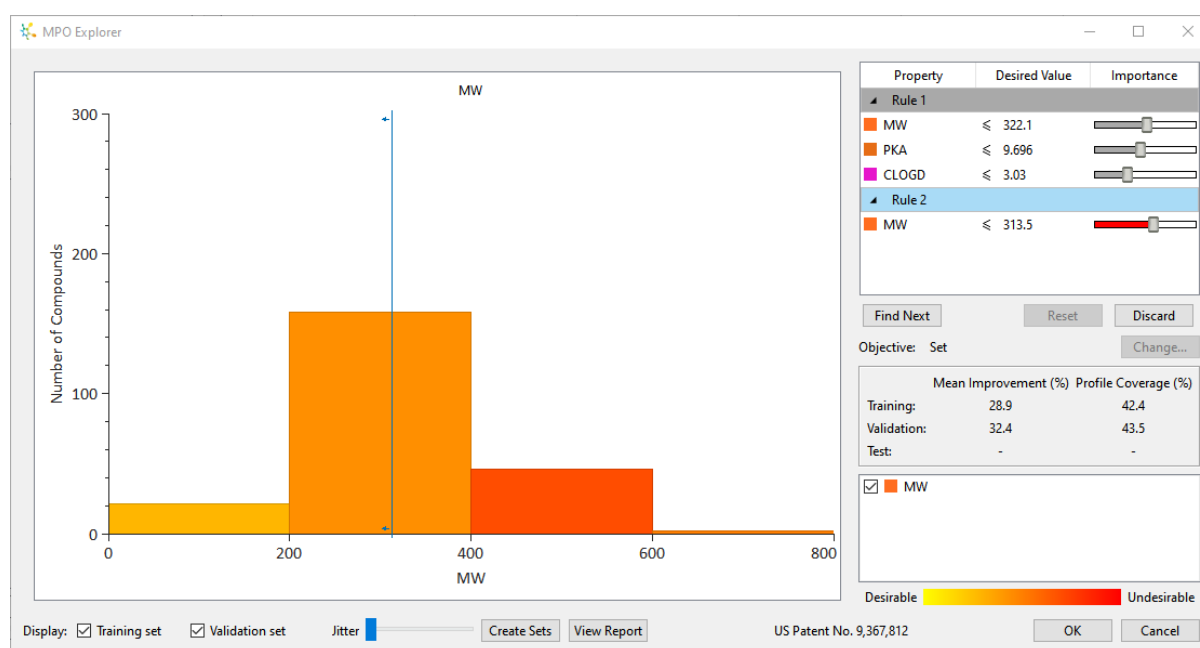
Note: the exact values you see will depend upon the precision you have specified for displaying numbers in the preferences.

Below this, the statistics for the corresponding multi-parameter rule can be seen; in this case, the mean improvement for compounds obeying the rule is > 54%, i.e. compounds that meet these

property criteria have a 54.1% greater chance of being a drug than a compound selected randomly from the set. If you hover your mouse over this statistic, a tool-tip will display additional information; in this case, the p-value for this rule is 0.0098, suggesting that it is highly statistically significant and the odds ratio is 4.6, which means that compounds meeting all three property criteria have a ~4× higher chance of success than compounds that do not. A detailed report on the statistics can be generated by clicking the **View Report** button.

On the left, plots show the rule in property space corresponding to the property criteria. The properties that are shown are controlled by the tick boxes in the bottom right of the analysis tool. In each plot, the blue lines indicate the boundaries implied by the rules. These boundaries can be dragged to modify the criteria, and the statistics will be updated instantly. The compounds in the training set are represented by circles and those in the validation set by 'x'. Desirable compounds (i.e. drugs) are shown in yellow and undesirable compounds (i.e. unsuccessful candidates) are shown in red. Grey points indicate compounds that have been filtered out by criteria other than those represented in the plot.

- In this case, we will accept the rule that has been generated automatically. To search for a second rule, click the **Find Next** button (as shown above).




The second rule, shown above, depends only on MW. In this case, the mean improvement is only 32.4%, and the corresponding p-value is 0.0443, suggesting that this rule is not as statistically significant as the first rule.

- Therefore, we will reject this rule by clicking the **Discard** button.

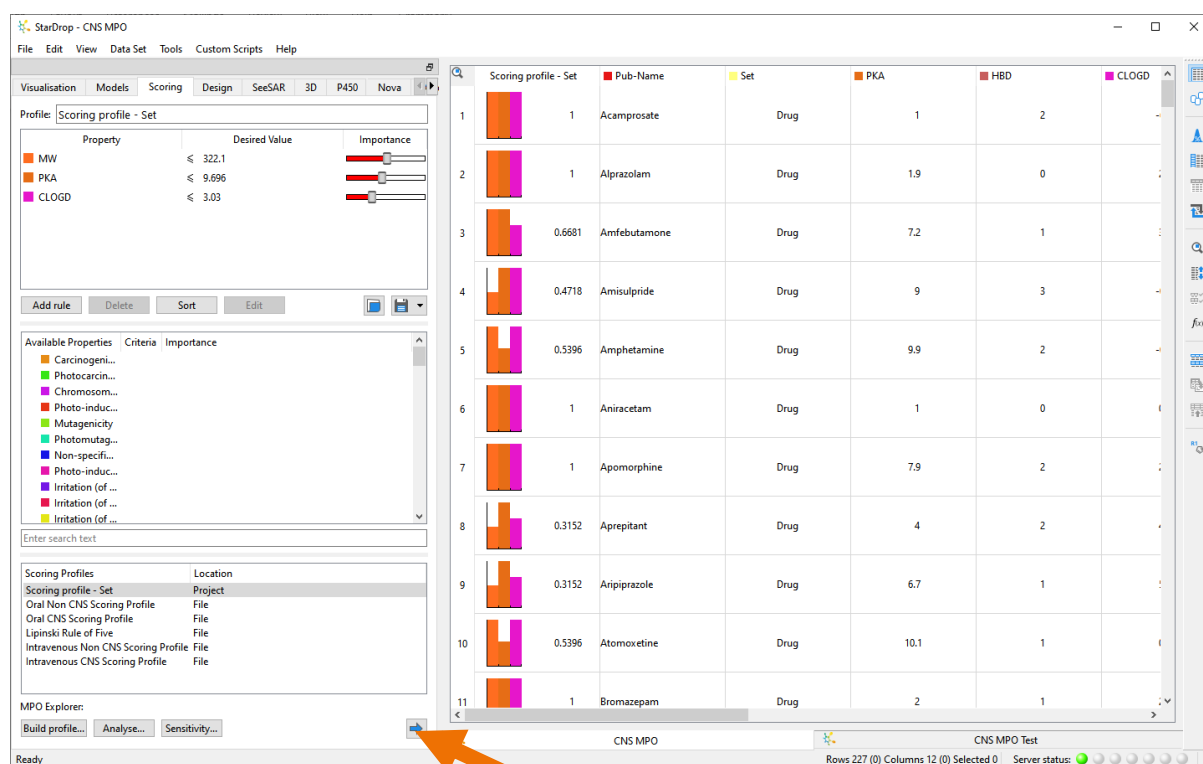
- The final profile, which in this case contains just one rule, will be displayed, and we can accept this by clicking the **OK** button.

The scoring profile will be displayed in the **Scoring** area. As with any scoring profile, we can modify, rename or save the profile. Clicking **Analyse** under **MPO Explorer** will return to the MPO Explorer profile analysis tool.

Note: you can analyse any scoring profile, not only those built using MPO Explorer's Profile Builder.

- Run the new scoring profile on the full data set by clicking the  button in the **Scoring** area.

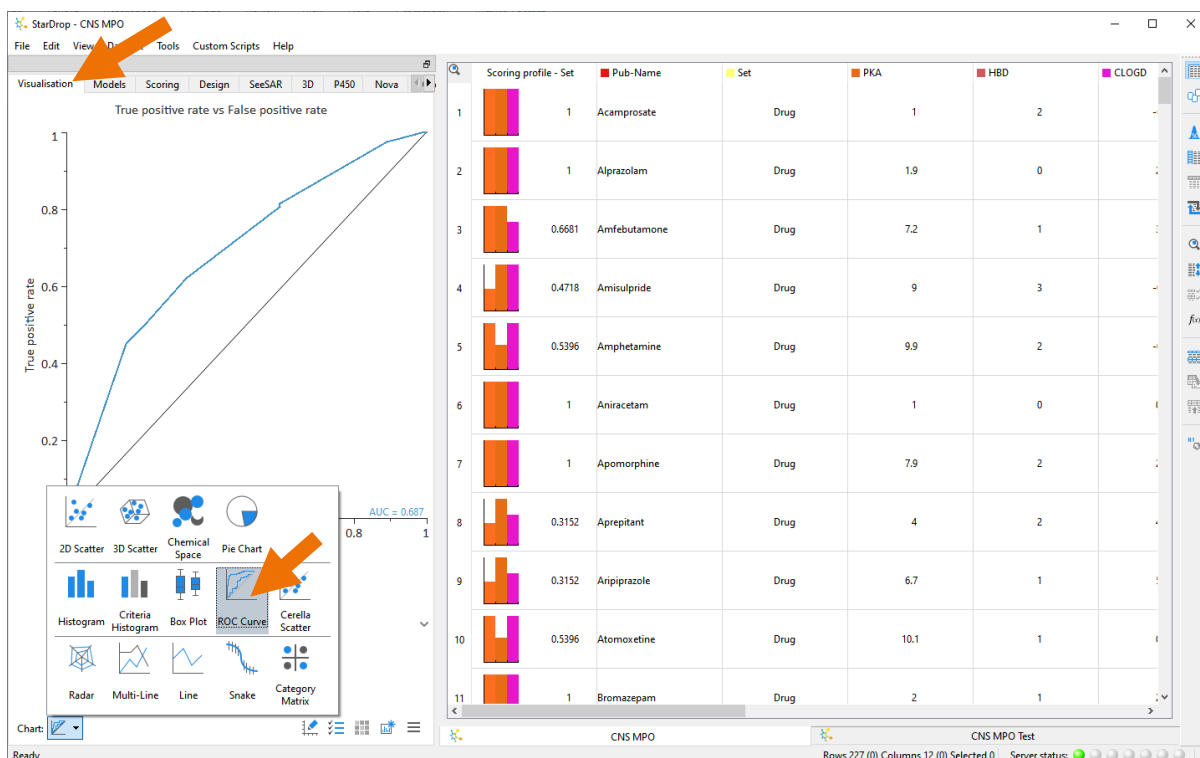
Note: you can ignore the warning about data with zero uncertainty.



We can compare the performance of this scoring profile with the CNS MPO score by plotting a Receiver Operating Characteristic (ROC) plot.


- Click on the **Visualisation** tab.
- From the **Chart** menu, select **ROC Curve**.

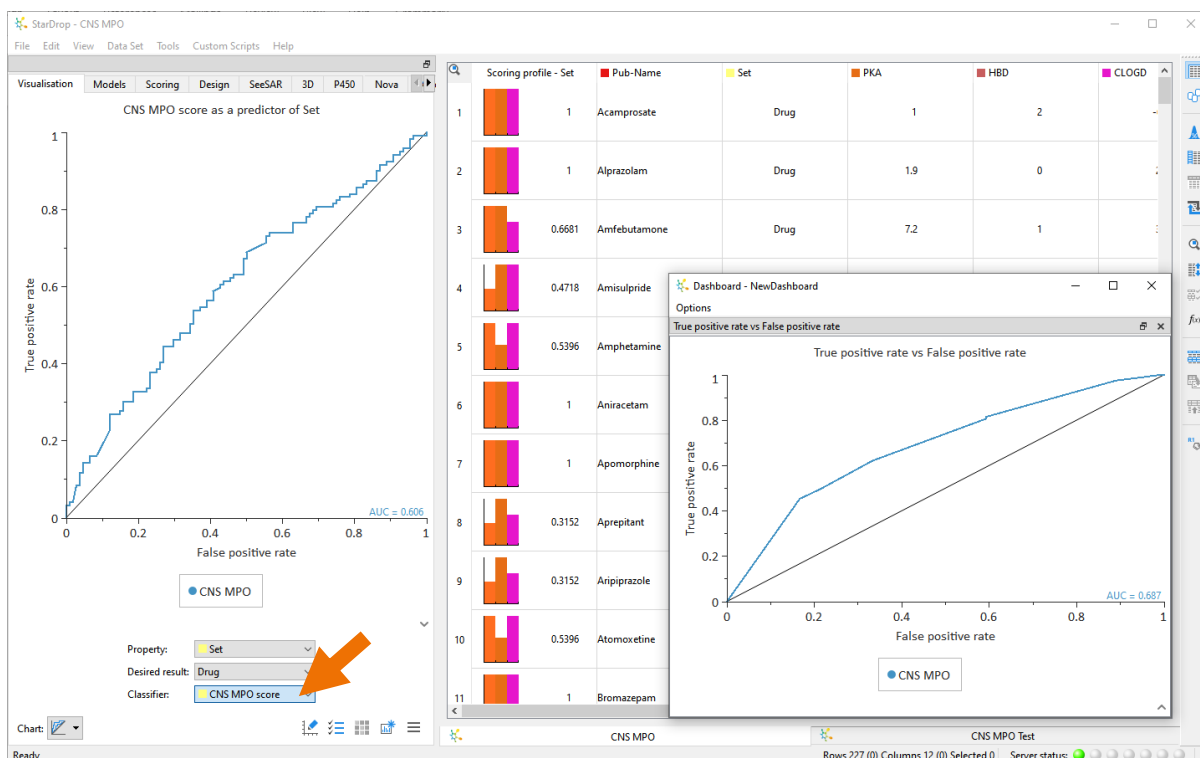
The chart will already have **Set** as the property, **Drug** as the desired result and use the new score **Scoring profile – Set** as the classifier.



We would like to see a ROC curve above the black identify line (which corresponds to the performance of a random selection) and ideally as close to the top-left corner as possible. A higher area under the curve (AUC) corresponds to better predictive performance. For more details of ROC curves, see http://en.wikipedia.org/wiki/Receiver_operating_characteristic.

Here we can see that the AUC for the scoring profile we have generated is approximately 0.69. We can compare this performance with that of the CNS MPO score.

- Click the **Detach** button  to add this ROC Curve to a new dashboard.
- In the **Visualisation** area, change the **Classifier** to **CNS MPO score** to create a ROC curve showing the performance of CNS MPO score, as shown below.




From this, we can see that the performance of the CNS MPO score is not as good as the scoring profile we have generated because the ROC curve is closer to the identity line, and the AUC value is lower.

This is not a fair test of CNS MPO score and the new profile that we have derived for two reasons:


- Compounds nominated as clinical candidates will generally have reasonable properties, so we would expect it to be quite challenging to distinguish candidates from successful drugs based on these simple properties. A more realistic test would be to distinguish drugs from early 'lead' compounds from drug discovery projects.
- We have assessed the performance of the scoring profile on the same set used to train and validate the corresponding rule. Therefore, the measure of performance may be artificially high.

To address these concerns, we can apply these scoring approaches to an independent test set.

The second data set in this project, called CNS MPO Test, contains 118 drugs (different from those used to find the rule) and 1000 compounds randomly selected from compounds in the ChEMBL database (<https://www.ebi.ac.uk/chembl/>) with a pK_i/pIC_{50} of greater than 6 (i.e. a K_i/IC_{50} less than 1 μM) against a CNS target. The target and measured pK_i/pIC_{50} of each compound from ChEMBL are included in the data set.

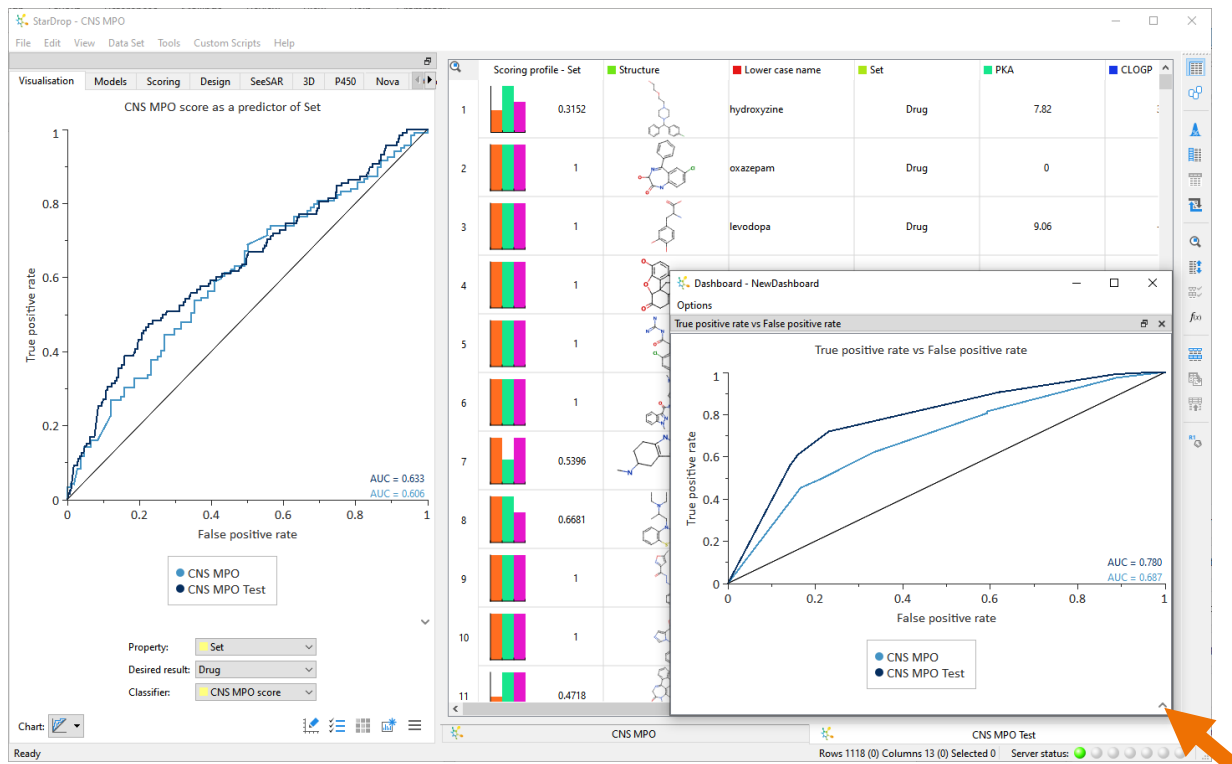
- Click on the tab at the bottom to display the second data set, change to the **Scoring** area and run the scoring profile by clicking the  button.

| Scoring profile - Set | Structure | Lower case name | Set | PKA | CLOGP |
|-----------------------|-----------|-----------------|------|------|-------|
| 1 | 0.3152 | hydroxyzine | Drug | 7.82 | |
| 2 | 1 | oxazepam | Drug | 0 | |
| 3 | 1 | levodopa | Drug | 9.06 | |
| 4 | 1 | | | | |
| 5 | 1 | | | | |
| 6 | 1 | | | | |
| 7 | 0.5396 | | | | |
| 8 | 0.6681 | | | | |
| 9 | 1 | | | | |
| 10 | 1 | | | | |
| 11 | 0.4718 | | | | |

- In the **Visualisation** area, click on the **Chart Data** button  at the bottom and tick the box next to CNS MPO Test to add this data set to the chart.

- In the dashboard, do the same to see how the new scoring profile performs against the independent test set.

Note: you may need to click on the small arrow in the bottom corner of the dashboard to display the controls.



From this, we can see that the performance of the scoring profile is better on this set, achieving an AUC of 0.78. However, AUC for the CNS MPO score has improved only marginally to 0.633, indicating that it cannot confidently distinguish between 'lead' compounds and drugs.

Conclusion

This example has shown how we can use the Profile Builder in MPO Explorer to generate scoring profiles with which to select compounds with a higher chance of success against our objective, in this case, distinguishing CNS drugs from unsuccessful candidates.

In this example, we have only used the simple functionality of the Profile Builder. Other capabilities enable the automatic selection of properties from a large number of possibilities and the derivation of 'soft' criteria to take into account the sparseness of data, helping to avoid 'hard' cut-offs that draw artificially harsh distinctions between compounds close to a property criterion.

It is notable that the simple scoring profile, using only three properties (logD, pK_a of the most basic site and MW), can outperform the CNS MPO score, which uses six properties. This illustrates the fact that there is a significant correlation between the properties used in CNS MPO score, for example,

logP and logD are strongly correlated ($R^2=0.6$ in this set). The inclusion of correlated properties can result in 'over-counting' of the same factor, inappropriately biasing the selection of compounds. The Profile Builder will select only property criteria that contribute significantly to the selection of high-quality compounds, avoiding the selection of multiple, highly correlated properties.

This example has used the simple 'drug-like' properties included in the CNS MPO score in order to draw a direct comparison. However, the Profile Builder can be applied to any data, including predicted or experimental biological or physicochemical properties that are more directly related to the *in vivo* disposition or efficacy required in a successful compound.

Additional, practical examples can be found in Yusof and Segall, Drug Discov. Today **19**(5) pp. 680-687, a preprint of which can be downloaded from:

<https://www.optibrium.com/publications-and-presentations/preprint-finding-the-rules-for-successful-drug-optimization/>.