

## Assignment 02

<b>Assignment Name:</b>	<b>Week 2 Assignment - Descriptor analysis of problem sets</b>
<b>Weight:</b>	6.25% of final grade
<b>Due Date:</b>	Tuesday, Midnight
<b>Associated Learning Outcomes:</b>	<p>6. Recognize and calculate basic physicochemical properties of small molecule compounds.</p> <p>7. Describe the underlying basis for cheminformatics techniques, and apply them to compare small molecule compounds. Applications for drug discovery.</p>
<b>Assignment Prompt:</b>	<p><b>Problem 1 - Data Analysis.</b> Download the fish toxicity dataset from the course, which contains 87 structures in smiles format and toxicity values for a particular fish. Higher LC50's indicate less toxicity. Investigate if there is any relationship between the toxicity observed and the following descriptors: molecular weight, number of rotatable bonds, number of aromatic bonds, ClogP (octanol/water partition coefficient), and TPSA (topological polar surface area). <u>Other descriptors should be explored and discussed if they are relevant.</u></p> <p><b>Hint:</b> LC50 is a continuous variable, but it may be more informative to “bin” continuous variables into categories such as “non-toxic vs toxic, or high, medium and low toxicity. Then you can graph the groups of compounds separately, and compare properties between groups.</p> <p>Turn in</p> <p>A) Your analysis and conclusions. <b>(6 points)</b> Are there any hypotheses you can make about the descriptors as to how they pertain to physical phenomenon? (i.e. are “greasy” compounds more toxic?)</p> <p>B) An image of your knime workflows and supporting graphs and graphics <b>(4 points)</b>, with captions.</p> <p>A complete analysis is likely to include box plots, and histograms. Use different color fills to illustrate toxic vs non-toxic so that we can see the differences between groups.</p>

**Problem 2 - Simple Model Building:** The blood-brain barrier (BBB) penetration ability of a compound is an important consideration in drug design. Drugs that act on targets in the central nervous system need to pass through the BBB, while a compound that acts on a peripheral target is generally safer if it does not pass the BBB. In this exercise, we'll look at how physiochemical properties and descriptors correlate with the BBB penetration ability of a compound.

Download the blood-brain barrier compound dataset on the course site. The BBB penetration ability is indicated as positive (p) or negative (n). Convert the smiles file to an SDF file using OpenBabel, either in Knime or separately. You can obtain OpenBabel from <http://openbabel.org>. Construct a report on this data that includes the following sections:

- 1) Exploratory data analysis (**5 points**)
  - a) Summarize the properties you feel relevant to achieving either penetrant (p) or nonpenetrant (n) compounds. Are there any relevant cutoffs that seem to materialize in the properties to discriminate?
  - b) Provide histograms or density plots relevant to the discussion, in a similar manner to above
  - c) Provide any relevant code or workflow screenshots as captioned figures
- 2) Simple Model Building (**10 points**)
  - a) Discuss 3 (or more) properties from above of your choice, and describe how they discriminate from penetrant or non-penetrant, and any caveats that you can think of. Pick a cutoff (decision boundary) for these 3 properties. Generate a simple model using one, two, three or more properties to separate permeable from non-permeable compounds (e.g. Property A > 4, Property B < 3, Property C > 0.5)
  - b) Run your model against the dataset (**Hint:** Use the rule engine node). Provide a Confusion Matrix for your model and analyze performance.
  - c) How can you make your model better? What happens to performance if you change the model? Are there any specific examples of false positives or false negatives that are instructive?

	<b>Note:</b> Please <b>do not</b> try to build flexible function (machine learning) models in this section. The idea here is to understand simple correlations first. Remember, the simple model is often the best model.
<b>Context/Purpose:</b>	In this exercise we introduce numerical interpretation of small molecules relative to a property we care about. The concept of model-building and evaluation is introduced, and we apply it to a practical example.
<b>Requirements and Logistics:</b>	Students should turn in a typed word or google docs document converted to PDF. All figures should be pasted into the document to generate a one-file submission. Handwritten materials will not be accepted.

## Assignment Rubric

	<b>Excellent</b>	<b>Acceptable</b>	<b>Incorrect / Absent</b>
<b>Assignment Questions</b>	Questions are answered and workflows screenshots are provided. <b>Full point value</b>	Partially correct, or supporting workflows are not provided. <b>Partial Point Value</b>	Incorrect, or not answered <b>No Point Value</b>

Late policy - **5 point penalty per day late - More than 3 days late: No credit**