

# How to predict new compounds

## (1) About the models of the prediction system:

**Table 1** shows the five models of the permeability prediction system.

**Table 1.** five models of the permeability prediction system.

Dataset	RRCK-C	PAMPA-C	Caco2-L	Caco2-C	Caco2-A
Model	MOE2D-SVR	MOE2D-LightGBM	MOE2D-RF	MOE2D-RF	MOE2D-RF

## (2) Descriptors :

The models were built on KNIME platform (version 4.7.0). The MOE2D descriptors were calculated by MOE (version 2019). The details of the selected MOE2D descriptors related to the above five models are as follows:

### (1) RRCK-C (103):

['apol', 'ast\_violation', 'a\_acc', 'a\_aro', 'a\_don', 'a\_donacc', 'a\_heavy', 'a\_hyd', 'a\_IC', 'a\_nC', 'a\_nF', 'a\_nH', 'a\_nN', 'a\_nO', 'a\_nS', 'balabanJ', 'BCUT\_PEOE\_1', 'BCUT\_PEOE\_2', 'BCUT\_SLOGP\_1', 'BCUT\_SLOGP\_2', 'BCUT\_SMR\_1', 'bpol', 'b\_1rotN', 'b\_double', 'b\_max1len', 'chi1v', 'chi1v\_C', 'chi1\_C', 'chiral', 'chiral\_u', 'diameter', 'h\_ema', 'h\_emd', 'h\_emd\_C', 'h\_logD', 'h\_logP', 'h\_logS', 'h\_log\_pbo', 'h\_pavgQ', 'h\_pKa', 'h\_pKb', 'h\_pstrain', 'Kier2', 'Kier3', 'KierA2', 'KierFlex', 'lip\_acc', 'lip\_violation', 'logP(o/w)', 'logS', 'opr\_brigid', 'opr\_nring', 'opr\_nrot', 'opr\_violation', 'PEOE\_PC+', 'PEOE\_PC-', 'PEOE\_VSA+0', 'PEOE\_VSA+1', 'PEOE\_VSA+2', 'PEOE\_VSA+3', 'PEOE\_VSA+4', 'PEOE\_VSA+5', 'PEOE\_VSA+6', 'PEOE\_VSA-0', 'PEOE\_VSA-1', 'PEOE\_VSA-2', 'PEOE\_VSA-3', 'PEOE\_VSA-4', 'PEOE\_VSA-5', 'PEOE\_VSA-6', 'PEOE\_VSA\_HYD', 'PEOE\_VSA\_NEG', 'PEOE\_VSA\_PNEG', 'PEOE\_VSA\_POL', 'PEOE\_VSA\_POS', 'PEOE\_VSA\_PPOS', 'petitjeanSC', 'radius', 'reactive', 'rsynth', 'SlogP', 'SlogP\_VSA0', 'SlogP\_VSA1', 'SlogP\_VSA2', 'SlogP\_VSA3', 'SlogP\_VSA4', 'SlogP\_VSA5', 'SlogP\_VSA6', 'SlogP\_VSA7', 'SlogP\_VSA8', 'SlogP\_VSA9', 'SMR\_VSA1', 'SMR\_VSA2', 'SMR\_VSA3', 'SMR\_VSA4', 'SMR\_VSA5', 'SMR\_VSA6', 'SMR\_VSA7', 'TPSA', 'VDistEq', 'vsa\_don', 'weinerPath', 'weinerPol']

### (2) PAMPA-C (71):

['apol', 'a\_aro', 'a\_don', 'a\_donacc', 'a\_hyd', 'a\_nO', 'balabanJ', 'BCUT\_SLOGP\_2', 'b\_1rotN', 'b\_max1len', 'chi1v\_C', 'chi1\_C', 'chiral', 'diameter', 'GCUT\_PEOE\_3', 'h\_emd', 'h\_emd\_C', 'h\_logD', 'h\_logP', 'h\_logS', 'h\_log\_pbo', 'h\_pavgQ', 'h\_pKa', 'h\_pKb', 'h\_pstates', 'h\_pstrain', 'Kier3', 'KierFlex', 'lip\_violation', 'logP(o/w)', 'logS', 'opr\_brigid', 'opr\_nring', 'opr\_nrot', 'opr\_violation', 'PEOE\_PC-', 'PEOE\_VSA+0', 'PEOE\_VSA+1', 'PEOE\_VSA+2', 'PEOE\_VSA+3', 'PEOE\_VSA+4', 'PEOE\_VSA+6', 'PEOE\_VSA-0', 'PEOE\_VSA-2', 'PEOE\_VSA-4', 'PEOE\_VSA-6', 'PEOE\_VSA\_NEG', 'PEOE\_VSA\_PNEG', 'PEOE\_VSA\_POS', 'petitjeanSC', 'radius', 'reactive', 'SlogP', 'SlogP\_VSA0', 'SlogP\_VSA1', 'SlogP\_VSA3', 'SlogP\_VSA4', 'SlogP\_VSA5', 'SlogP\_VSA8', 'SlogP\_VSA9', 'SMR\_VSA1', 'SMR\_VSA2', 'SMR\_VSA3', 'SMR\_VSA4', 'SMR\_VSA5', 'SMR\_VSA6', 'SMR\_VSA7', 'VAdjMa', 'VDistEq', 'vsa\_don', 'weinerPath']

### (3) Caco2-L (82):

['apol', 'ast\_fraglike', 'ast\_violation', 'ast\_violation\_ext', 'a\_acc', 'a\_aro', 'a\_don', 'a\_hyd', 'a\_nN', 'a\_nS', 'balabanJ', 'BCUT\_PEOE\_2', 'BCUT\_SLOGP\_2', 'BCUT\_SLOGP\_3', 'BCUT\_SMR\_2', 'b\_double', 'b\_max1len', 'b\_rotR', 'chiral', 'chiral\_u', 'diameter', 'GCUT\_PEOE\_3', 'GCUT\_SLOGP\_3', 'h\_emd\_C', 'h\_logD', 'h\_logP', 'h\_logS', 'h\_log\_pbo', 'h\_pavgQ', 'h\_pKa', 'h\_pKb', 'h\_pstates', 'h\_pstrain', 'lip\_don', 'lip\_druglike', 'lip\_violation', 'logP(o/w)', 'logS', 'mutagenic', 'opr\_brigid', 'opr\_leadlike', 'opr\_violation', 'PEOE\_PC-', 'PEOE\_VSA+0', 'PEOE\_VSA+1', 'PEOE\_VSA+2', 'PEOE\_VSA+3', 'PEOE\_VSA+4', 'PEOE\_VSA+5', 'PEOE\_VSA+6', 'PEOE\_VSA-0', 'PEOE\_VSA-1', 'PEOE\_VSA-2', 'PEOE\_VSA-3', 'PEOE\_VSA-4', 'PEOE\_VSA-6', 'PEOE\_VSA\_FNEG', 'PEOE\_VSA\_FPOS', 'PEOE\_VSA\_NEG', 'radius', 'reactive', 'rsynth', 'SlogP', 'SlogP\_VSA1', 'SlogP\_VSA3', 'SlogP\_VSA4', 'SlogP\_VSA5', 'SlogP\_VSA6', 'SlogP\_VSA7', 'SlogP\_VSA8', 'SlogP\_VSA9', 'SMR\_VSA1', 'SMR\_VSA3', 'SMR\_VSA4', 'SMR\_VSA5', 'SMR\_VSA6', 'SMR\_VSA7', 'VAdjEq', 'VAdjMa', 'VDistEq', 'vsa\_don', 'weinerPath']

#### **(4) Caco2-C (79):**

['apol', 'ast\_violation', 'ast\_violation\_ext', 'a\_acc', 'a\_aro', 'a\_don', 'a\_donacc', 'a\_nCl', 'a\_nF', 'a\_nN', 'a\_nS', 'balabanJ', 'b\_1rotN', 'b\_double', 'b\_max1len', 'chiral', 'chiral\_u', 'diameter', 'GCUT\_PEOE\_3', 'h\_emd', 'h\_emd\_C', 'h\_logD', 'h\_logP', 'h\_logS', 'h\_log\_dbo', 'h\_log\_pbo', 'h\_pavgQ', 'h\_pKa', 'h\_pKb', 'h\_pstates', 'h\_pstrain', 'lip\_druglike', 'lip\_violation', 'logP(o/w)', 'logS', 'opr\_brigid', 'opr\_leadlike', 'opr\_nring', 'opr\_violation', 'PEOE\_PC-', 'PEOE\_VSA+0', 'PEOE\_VSA+1', 'PEOE\_VSA+2', 'PEOE\_VSA+3', 'PEOE\_VSA+4', 'PEOE\_VSA+5', 'PEOE\_VSA+6', 'PEOE\_VSA-0', 'PEOE\_VSA-1', 'PEOE\_VSA-2', 'PEOE\_VSA-3', 'PEOE\_VSA-4', 'PEOE\_VSA-6', 'PEOE\_VSA\_NEG', 'petitjeanSC', 'radius', 'reactive', 'SlogP', 'SlogP\_VSA0', 'SlogP\_VSA1', 'SlogP\_VSA3', 'SlogP\_VSA4', 'SlogP\_VSA5', 'SlogP\_VSA6', 'SlogP\_VSA7', 'SlogP\_VSA8', 'SlogP\_VSA9', 'SMR\_VSA1', 'SMR\_VSA2', 'SMR\_VSA3', 'SMR\_VSA4', 'SMR\_VSA5', 'SMR\_VSA6', 'SMR\_VSA7', 'TPSA', 'VAdjMa', 'VDistEq', 'vsa\_don', 'weinerPath']

#### **(5) Caco2-A (92):**

['apol', 'ast\_fraglike', 'ast\_violation', 'ast\_violation\_ext', 'a\_acc', 'a\_aro', 'a\_don', 'a\_donacc', 'a\_hyd', 'a\_nCl', 'a\_nF', 'a\_nN', 'a\_nO', 'a\_nS', 'balabanJ', 'BCUT\_SLOGP\_2', 'b\_1rotN', 'b\_1rotR', 'b\_double', 'b\_max1len', 'b\_rotR', 'chiral', 'chiral\_u', 'diameter', 'GCUT\_PEOE\_3', 'h\_logD', 'h\_logP', 'h\_logS', 'h\_log\_dbo', 'h\_log\_pbo', 'h\_pavgQ', 'h\_pKa', 'h\_pKb', 'h\_pstates', 'h\_pstrain', 'Kier3', 'lip\_don', 'lip\_druglike', 'lip\_violation', 'logP(o/w)', 'logS', 'opr\_brigid', 'opr\_leadlike', 'opr\_nring', 'opr\_violation', 'PEOE\_PC+', 'PEOE\_PC-', 'PEOE\_VSA+0', 'PEOE\_VSA+1', 'PEOE\_VSA+2', 'PEOE\_VSA+3', 'PEOE\_VSA+4', 'PEOE\_VSA+5', 'PEOE\_VSA+6', 'PEOE\_VSA-0', 'PEOE\_VSA-1', 'PEOE\_VSA-2', 'PEOE\_VSA-3', 'PEOE\_VSA-4', 'PEOE\_VSA-6', 'PEOE\_VSA\_FNEG', 'PEOE\_VSA\_FPOS', 'PEOE\_VSA\_NEG', 'petitjeanSC', 'radius', 'reactive', 'rings', 'rsynth', 'SlogP', 'SlogP\_VSA0', 'SlogP\_VSA1', 'SlogP\_VSA3', 'SlogP\_VSA4', 'SlogP\_VSA5', 'SlogP\_VSA6', 'SlogP\_VSA7', 'SlogP\_VSA8', 'SlogP\_VSA9', 'SMR\_VSA1', 'SMR\_VSA2', 'SMR\_VSA3', 'SMR\_VSA4', 'SMR\_VSA5', 'SMR\_VSA6', 'SMR\_VSA7', 'TPSA', 'VAdjEq', 'VAdjMa', 'VDistEq', 'vsa\_acc', 'vsa\_don', 'weinerPath']

### ***(3) How to form your own prediction pipeline:***

The details of constructing your own workflow are shown in **Figure 1**, **Figure 2** and **Figure 3**.

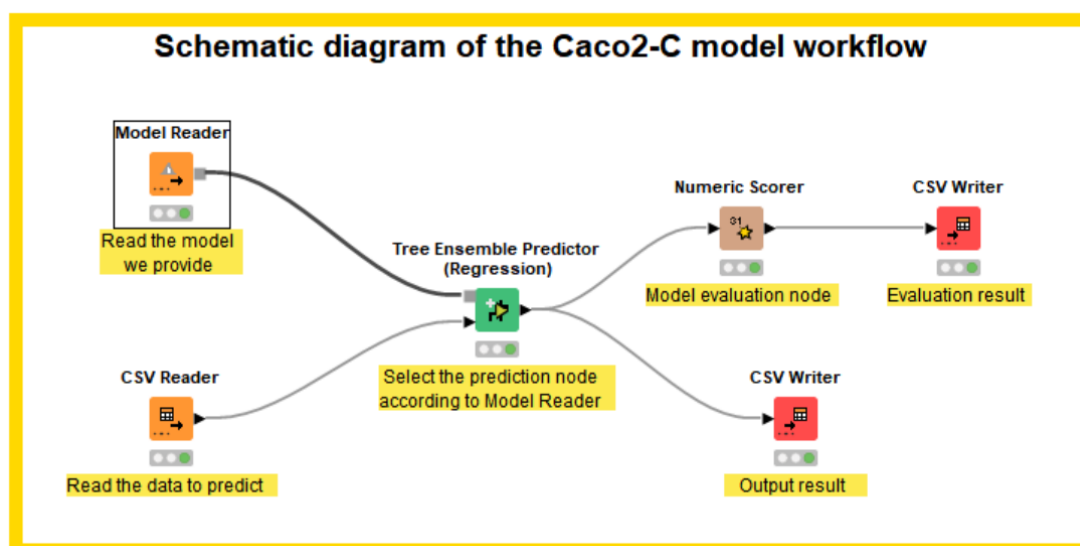
Explanation of workflow:

1) The local model file is read by the **Model Reader** node above, while the below reads the model-Normalizer.zip for normalizer;

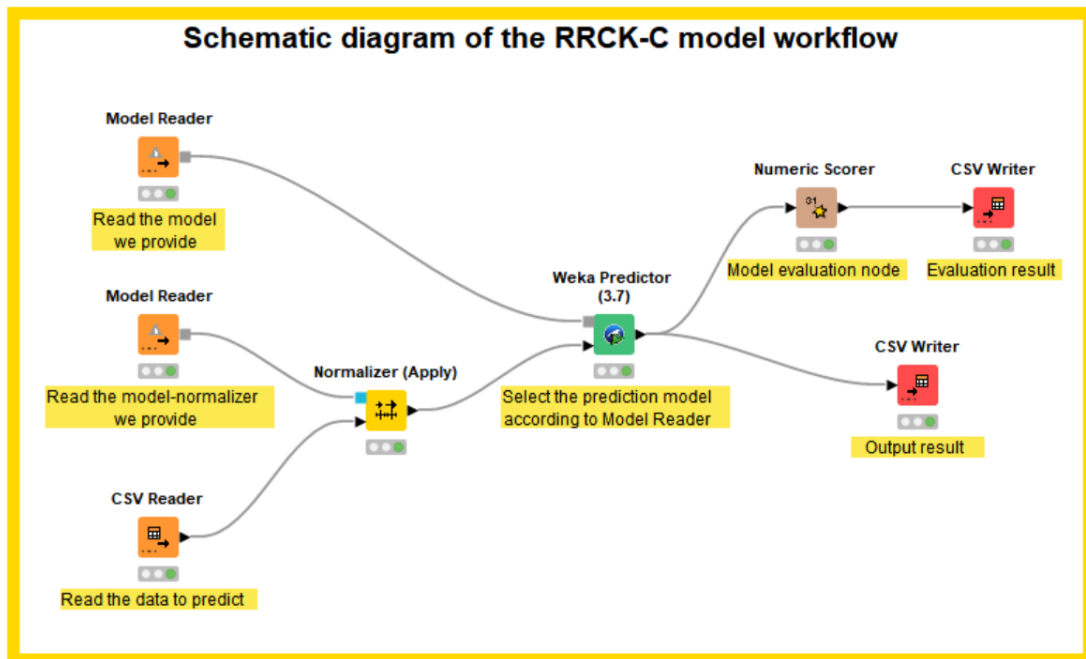
2) The node of **File Reader** is used to read the data that needs to be predicted; Please read your data as the examples we provided (example for RRCK-C.csv, example for Caco2-C.csv and example for PAMPA-C.csv). **If you just want to predict molecules without experimental values or labels, you should ignore the 'Permeability' and disconnect the evaluation nodes for example 'Numeric Scorer'.**

3) Select the corresponding prediction node according to the model read by the model reader;

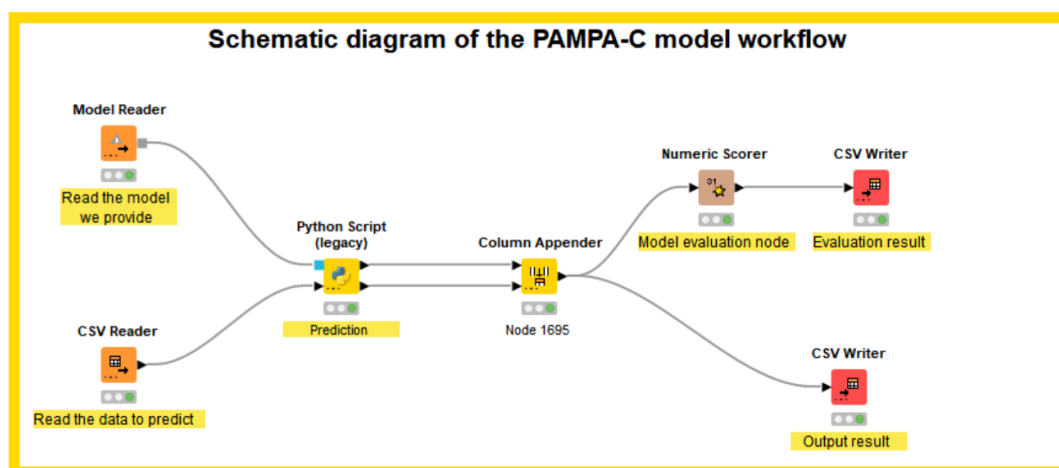
4) Output for the prediction result. In addition, evaluation nodes can be chosen according to your task.



**Figure 1.** KNIME usage example of Caco2-C model.



**Figure 2.** KNIME usage example of RRCK-C regression model.



**Figure 3.** KNIME usage example of PAMPA-C regression model.