# log *P* Predictions Using SM*x* or LSER With Training Data from DrugBank.ca

## Jonathan A. Ouimet, Rachel C. Ollier, and Andrew S. Paluch

Department of Chemical, Paper, and Biomedical Engineering
Miami University

PaluchAS@MiamiOH.edu

SAMPL6 Virtual Workshop

May 16, 2019

# MIAMI UNIVERSITY
OXFORD, OH · EST. 1809

# Our Entries

- SM12-Solvation-Trained (7/91)
- SM12-Solvation
- SM8-Solvation-Trained
- SM8-Solvation
- SMD-Solvation-Trained
- SMD-Solvation

- GC-LSER
- ISIDA-LSER
- UFZ-LSER

# The Idea

The standard approach:

$$\log P = -\frac{1}{\ln(10)}\left(\frac{\Delta G_{\text{oct}}^{\text{solv}}}{RT} - \frac{\Delta G_{\text{water}}^{\text{solv}}}{RT}\right)$$
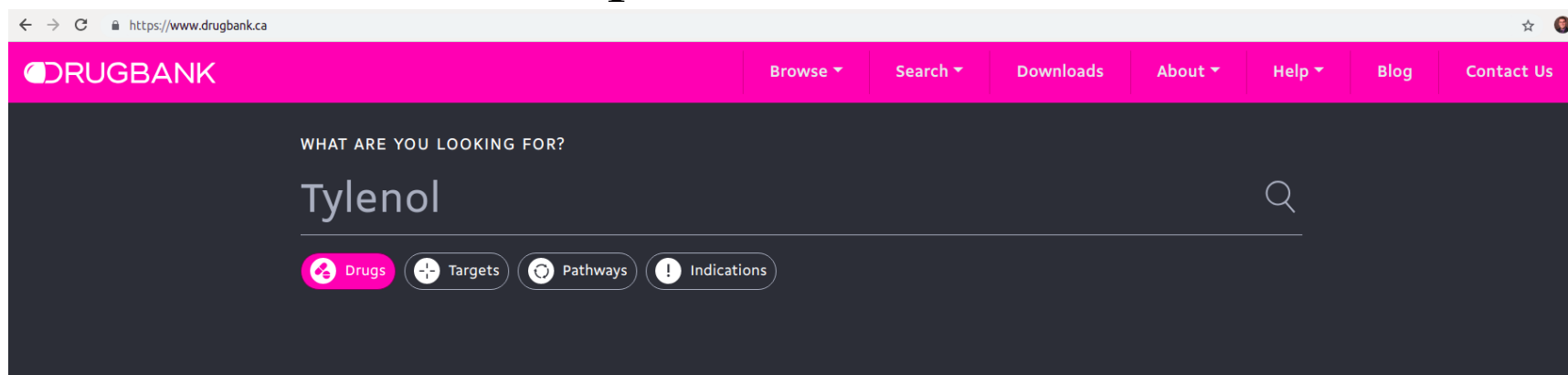
Which assumes:

1) Solute is infinitely dilute (i.e., no solute-solute interactions)
2) Water and octanol are immiscible

What if instead we used:

$$\log P = -a\frac{\Delta G_{\text{oct}}^{\text{solv}}}{RT} + b\frac{\Delta G_{\text{water}}^{\text{solv}}}{RT} + c$$

# Training Set

- We sought a training set representative of the SAMPL6 molecules for which predictions were to be made.

# Training Set

- We sought a training set representative of the SAMPL6 molecules for which predictions were to be made.

  - 6 of the molecules had the 4-amino quinzaoline scaffold

  - all the molecules had multiple aromatic rings within their structure

- Selected molecules that

  - Had the 4-amino quinzaoline scaffold

  - Experimental log $P$ values were available

  - size and molecular weight were considered to ensure that the selected molecules were representative of the SAMPL6 molecules

  - Include halogens and electronegative elements

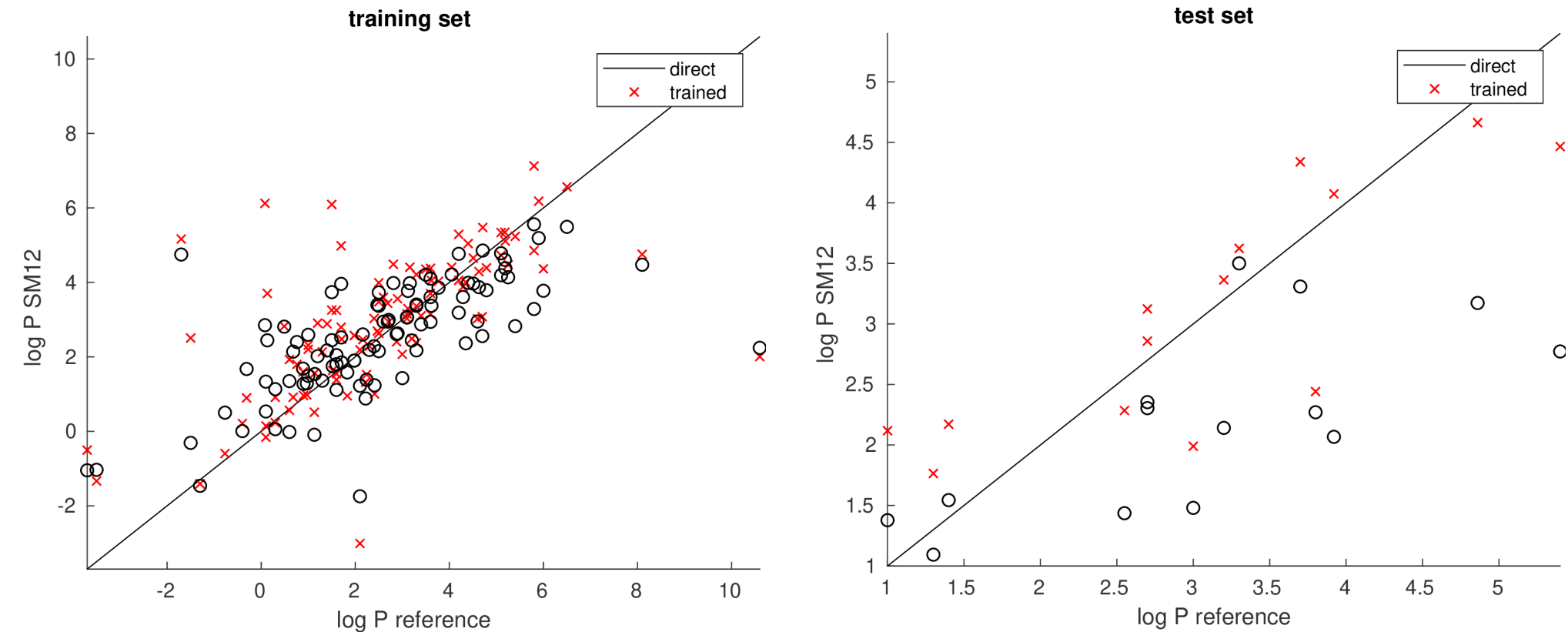- Final training set of 100 molecules from DrugBank.ca

# Test/Validation Set

- Performed a similarity in the DrugBank.ca database
  - Similarity score greater than 0.5 with the SAMPL6 structures
  - Also included "approved" drugs with know log $P$ values
- Test set of 14 molecules
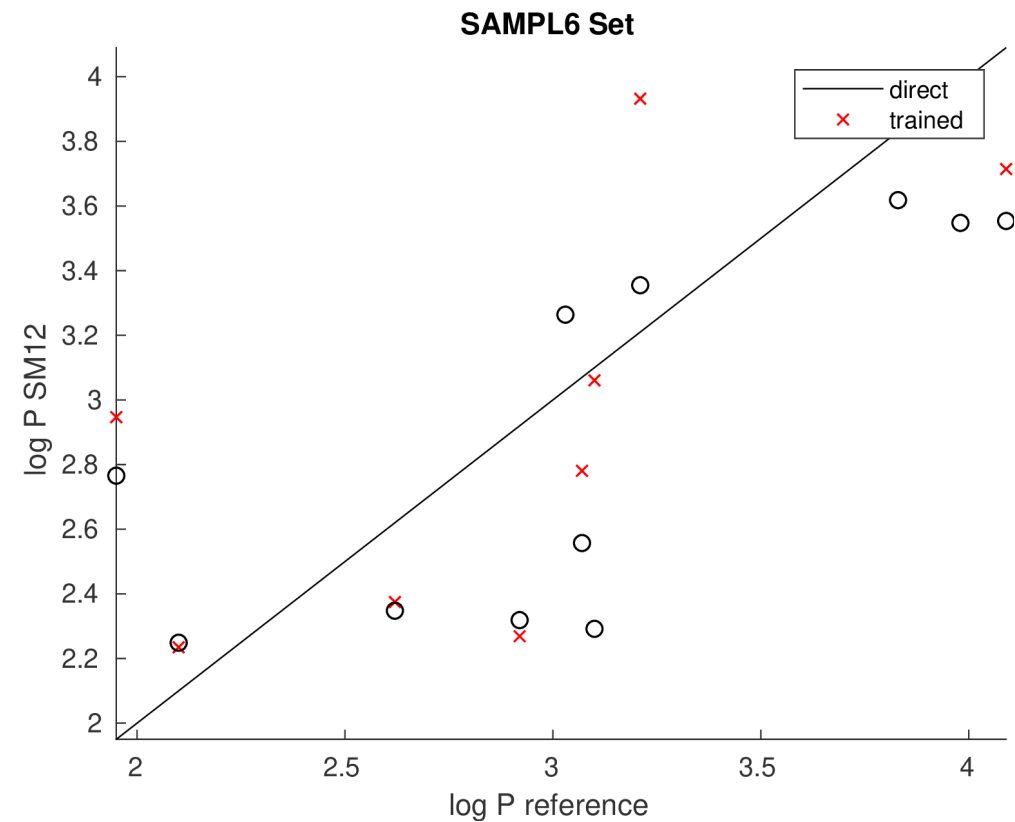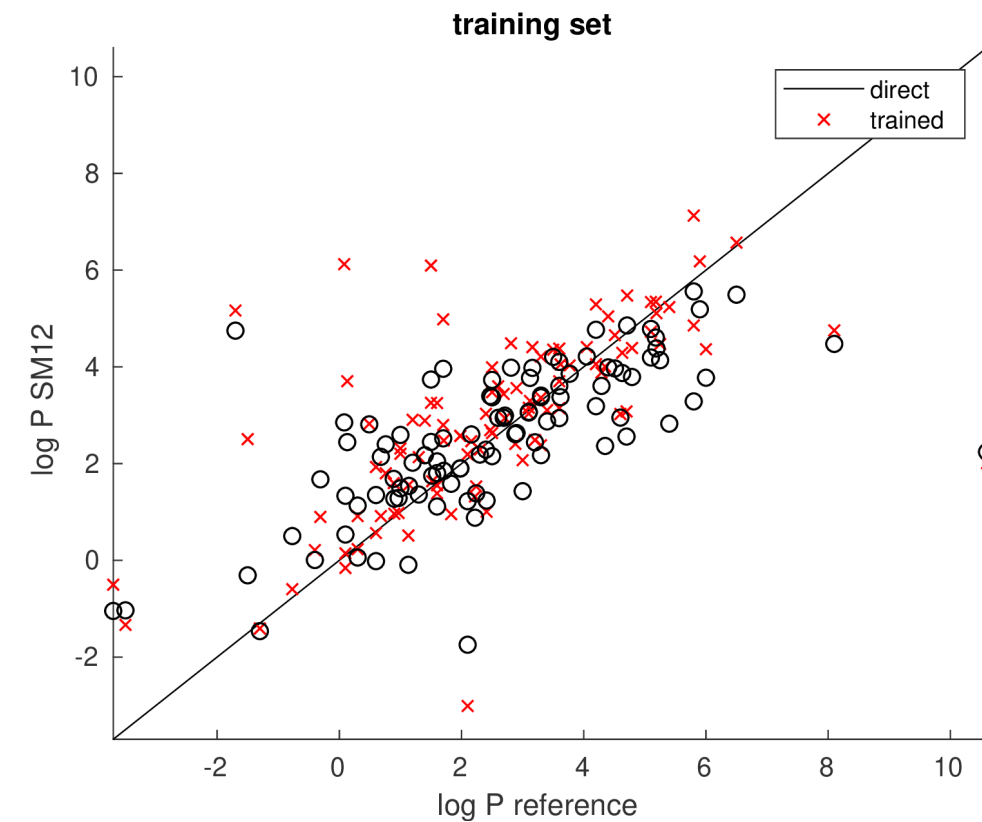- Used to estimate model accuracy

# Calculations

- Start with SMILES from DrugBank.ca and SAMPL6

- Open Babel using GAFF

  – Generate 3-D structures and then performed systematic conformation search to find lowest energy conformer

- QChem 5.1.2

  – Geometry optimization at the M06-2X/cc-pVDZ level of theory/basis set

  – Single point energy energy calculations in the SM12, SM8, and SMD continuum solvent models for 1-octanol and water at the M06-2X/6-31G(d) level of theory/basis set
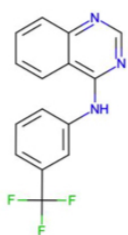
# Top Performing SM12



training set

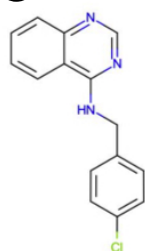test set

# Top Performing SM12
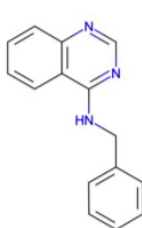
# Possible sources of error within our control

- Conformation

- Theory/basis set

- More descriptors and/or machine learning in place of multi-linear regression

- How to treat HCl

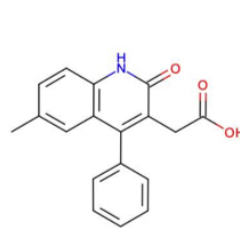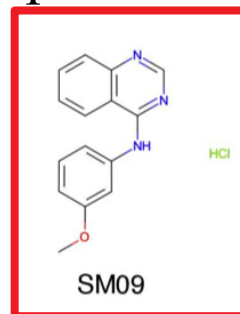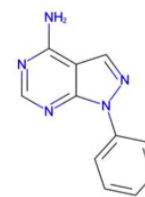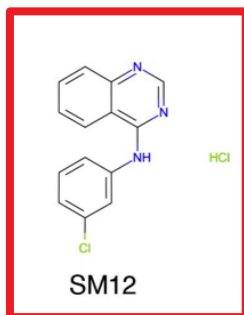  – Performing calculations without HCl for publication
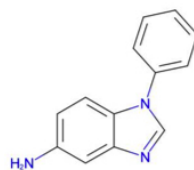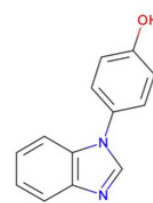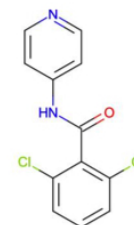
# Some Lessons Learned

- Computed solvation free energies are sensitive to choice of SM12, SM8, or SMD.

  - The results are expected to be sensitive to theory/basis set too.

  - And geometry.

- Including "training" data did effect the results. But it is unclear how significant this is.

  - Could also investigate choice of training data further.

  - DrugBank.ca was chosen due to its size and it is freely available.

- HCl

# Acknowledgements

- Ohio Supercomputer Center

PaluchAS@MiamiOH.edu