University of Heidelberg



Validation and Predictivity of QSAR Models

Hugo Kubinyi

University of Heidelberg, Germany

E-Mail kubinyi@t-online.de http://home.t-online.de/ home/kubinyi

15th EURO-QSAR Symposium Istanbul, Turkey, Sept. 2004



University of Heidelberg

A Common Situation

A chemist synthesizes about 30 compounds.

The biologist determines the activity values.

Both ask the chemoinformatician to derive a QSAR model.

The chemoinformatician loads 1500 variables (e.g. from the program DRAGON, Roberto Todeschini) and derives a QSAR model, containing only a few variables, which meets all statistical criteria.

Chemist, biologist and chemoinformatician publish the results. Everybody is happy.

| The Selwood Data Set | |
|--------------------------------|---------------------------------------|
| n = 31 compound | s and k = 53 independent variables. |
| Theoretically, the | re are: |
| 53 | one-variable models, |
| 1,378 | two-variable models, |
| 23,426 | three-variable models, |
| 292,825 | four-variable models, |
| , 22.957.480 | six-variable models. |
| | in total |
| 7,160,260,814 | ,092,303 regression models, |
| containing on selected from | e to 29 variables, 53 X-variables. |



The 53 X Variables of the Selwood Data Set

ATCH1 - ATCH10 = partial atomic charges DIPV_X, DIPV_Y and DIPV_Z = dipole vectors DIPMOM = dipole moment ESDL1 - ESDL10 = electrophilic superdelocalizability NSDL1 - NSDL10 = nucleophilic superdelocalizability VDWVOL = van der Waals volume SURF_A = surface area MOFI_X, MOFI_Y and MOFI_Z = moments of inertia PEAX_X, PEAX_Y and PEAX_Z = ellipsoid axes MOL_WT = molecular weight S8_1DX, S8_1DY and S8_1DZ = substituent dimensions S8_1CX, S8_1CY and S8_1CZ = substituent centers LOGP = partition coefficient M_PNT = melting point SUM_F and SUM_R = sums of the F and R constants



University of Heidelberg

Questions

Can we derive "good" (statistically valid) models: yes

Do our models have internal predictivity (Q² values): yes

Are these models "better" than models from scrambled or random data (y, x, y and x): yes

Are 53 X variables too many to select from: no, fine

Can our models predict a test set (r²_{pred} value): not at all

Is there a relationship between internal and external predictivity: by no means





























S. H. Unger and C. Hansch J. Med. Chem. <u>16</u>, 745-749 (1973)

One must rely heavily on statistics in formulating a quantitative model but, at each critical step in constructing the model, one must set aside statistics and ask questions. ... without a qualitative

perspective one is apt to generate statistical unicorns, beasts that exist on paper but not in reality.

... it has recently become all too clear that one can correlate a set of dependent variables using random numbers as dependent variables. Such correlations meet the usual criteria of high significance ...





University of Heidelberg Summary, Conclusions and Recommendations Apply the Unger and Hansch recommendations: 1. Selection of meaningful variables 2. Elimination of interrelated variables 3. Justification of variable choices by statistics 4. Principle of parsimony (Ockham's Razor) 5. Number of variables to choose from 6. Number of variables in the model 7. Qualitative biophysical model Additional recommendations: 8. Beware of Q² (Alex Tropsha) 9. Search for outliers in the test set 10. Do not expect your model to be predictive

