

A combined QSAR model for mutagenicity

T. Ferrari, G. Gini A. Roncaglioni, E. Benfenati
DEI, Politecnico di Milano, Italy Istituto Mario Negri, Milan, Italy
email: tferrari@elet.polimi.it

INTRODUCTION

Mutagenicity has to be evaluated for each chemical within the REACH legislation. Several experimental methods exist, and also *in silico* methods. Some *in silico* models use human knowledge, which has been codified into rules (expert systems), while other models extract the information using data mining approaches.

Within the CAESAR project we developed a QSAR model for mutagenicity, with the purpose to produce a model suitable for REACH.

The Algorithms

To achieve a model more suitable for REACH, two different cascading techniques have been arranged: a machine learning algorithm, to build an early model with the best statistical accuracy, then an expert system to refine its predictions.

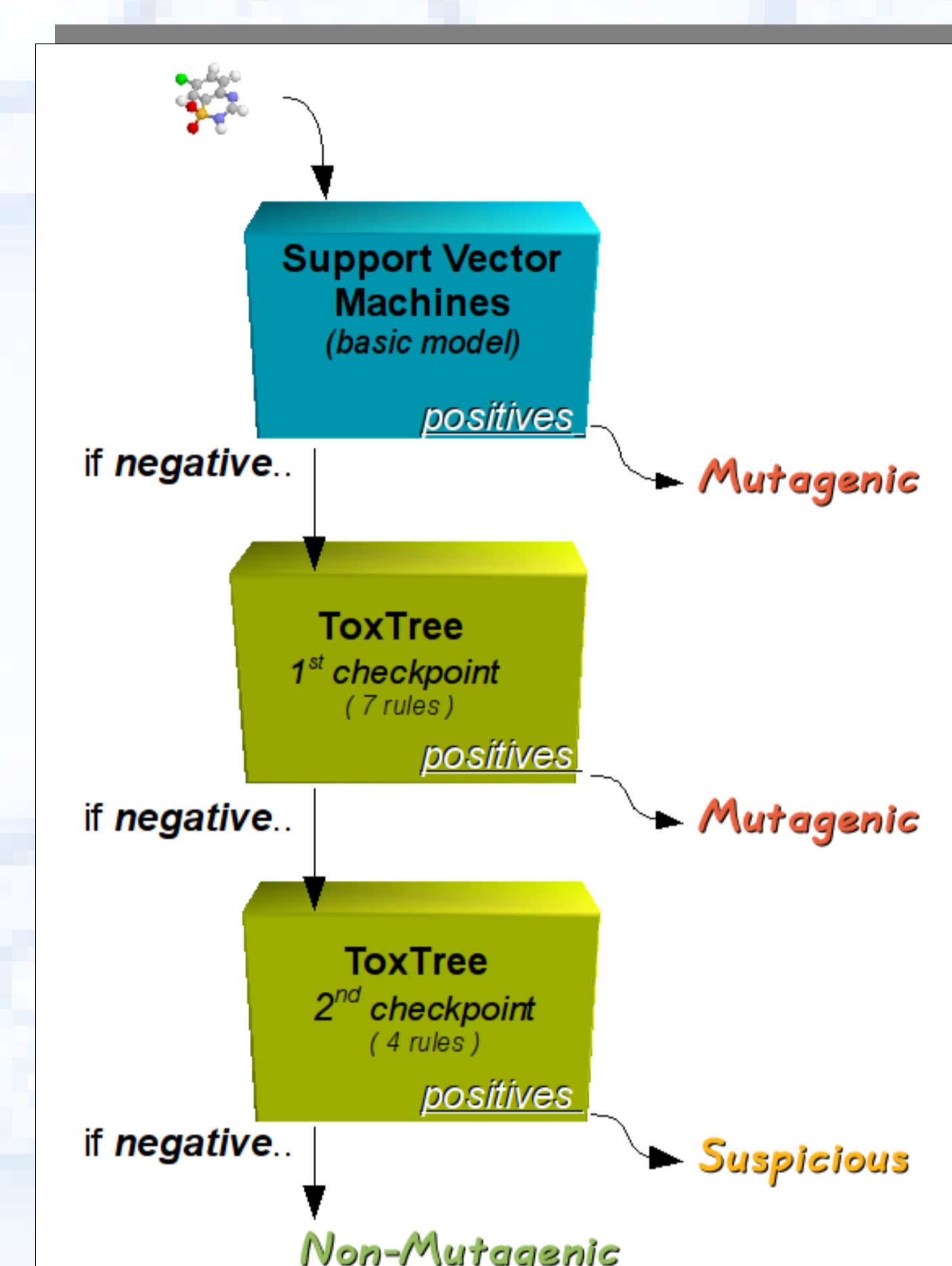
In the machine learning starting phase, a model was practiced on the training set using **LibSVM**⁴, an open source java implementation of *Support Vector Machines* (SVM) algorithm, with *Radial Basis Function* (RBF) kernel. The optimal parameterization for the kernel function was automatically performed by a grid-search in the parameter space. This provided a basic model with very good performances.

Then, a further scan for structural alerts for genotoxic carcinogenicity on compounds predicted* as negative (*non-mutagenic*) by SVM model was carried out by **ToxTree**. An analysis on joined results permitted to identify two different subsets of ToxTree rules that better teams up with SVM predictions on two purposes: improving the accuracy and avoiding *false negatives*.

The final model, integration between the two approaches, can be designed like in figure.

(*) 10-folds cross-validation on the training set

Figure. Compounds evaluated as positive by SVM (■) are immediately classified as **Mutagenic**, the presumed negative are further tested by two consecutive checkpoints based on structural alerts (■). The first has the aim to enhance the prediction accuracy attempting a precise identification of misclassified FN, the second goes on with the FN removal as much as this doesn't noticeably downgrade the original accuracy (generating too many FP as well). To point out this distinction, compounds picked out by the former checkpoint are classified as **Mutagenic**, and those picked out by the latter one are classified as **Suspicious**. Unaffected compounds are finally classified as **Non-Mutagenic**.



MATERIALS AND METHODS

In order to produce models for regulatory purposes it is important to address:

- data quality control
- reduction of *false negatives*
- model reproducibility
- model validation

The Data

A large data set of compounds (more than 4000) was used -as described by Katzius, McGuire and Bursi¹- and each chemical structure was individually checked within the CAESAR project to increase the model robustness.

For each compound, 37 molecular descriptors were calculated using **MDL QSAR** software²: 4 of them are global descriptors and the others are small 2D substructures counts. The resulting data were finally normalized in the [-1,1] range.

ToxTree³ was used on the same data set.

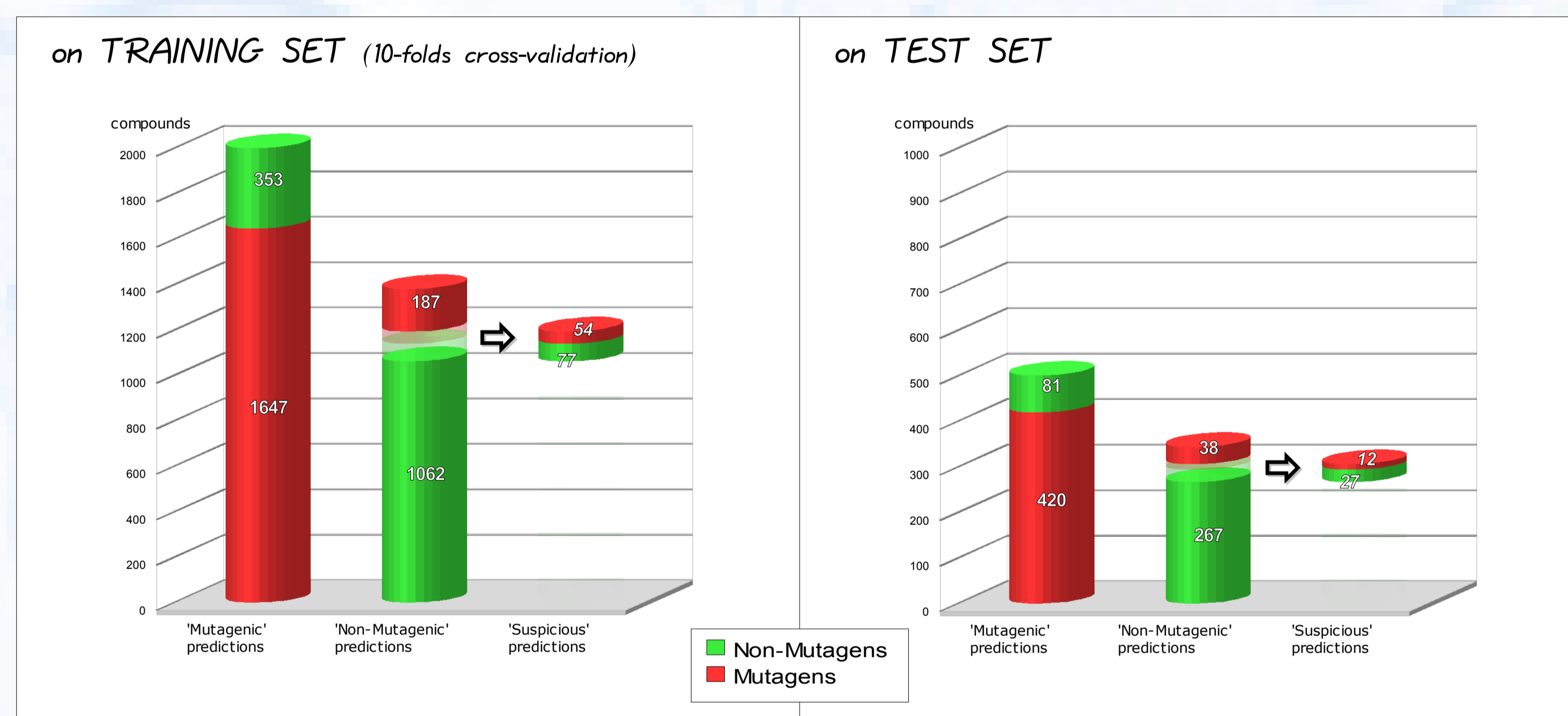
To provide a good basis for validation, the model was developed using a training set of about 80% of the compounds and the other compounds were left for testing.

RESULTS AND DISCUSSIONS

The combined model has three possible outputs: *Mutagenic*, *Non-Mutagenic* or *Suspicious*. The *Suspicious* prediction it's a warning: it just denotes a candidate mutagen, since it has fired a structural alert but not a so specific one. The choice for the final binary classification is left to the end-user, depending on his/her purposes. This leads to two different statistic results with different features: best accuracy or low false negative rate.

Results of combined model validation are presented in the table below.

VALIDATION OF THE MODEL :



ERROR-SHIFTING ANALYSIS :



CHOICE by end-user

Aimed at performances : OR Aimed at prudence :

↳ Suspicious compounds taken as Non-Mutagenic ↳ Suspicious compounds taken as Mutagenic

Validation on the TRAINING set:	Validation on the TEST set:	Validation on the TRAINING set:	Validation on the TEST set:
SVM in 10-folds cross-validation + 7 ToxTree rules	SVM + 7 ToxTree rules	SVM in 10-folds cross-validation + 11 ToxTree rules	SVM + 11 ToxTree rules
accuracy: 82.43%	accuracy: 84.50%	accuracy: 81.75%	accuracy: 82.72%
sensitivity: 87%	sensitivity: 89%	sensitivity: 90%	sensitivity: 92%
specificity: 76%	specificity: 78%	specificity: 71%	specificity: 71%

Combined model has good accuracy and a low *false negative* rate. This effect is powered by the 'expert' *ToxTree* layer supervision on compounds at first predicted as *Non-Mutagenic* by the former *SVM* layer, as explicated in the table on the right.

CAESAR SVM basic model has higher performances than ToxTree only. However, when used as a further screening tool, ToxTree was useful to avoid a certain number of *false negatives*. This shows the utility of integrating models based on different approaches: SVM is based on data mining and chemical descriptors, while ToxTree is based on knowledge from human experts and chemical fragments.

We remember that the reproducibility of the experimental results is about 85%.

REFERENCES

- (1) J. Med. Chem. **2005**
- (2) <http://www.mdl.com/>
- (3) From from ECB, JRC. Software available at <http://ecb.jrc.ec.europa.eu/qsar/>
- (4) Chang, CC. and Lin, CJ. LibSVM: a library for support vector machines, **2001**. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.

CONCLUSIONS

CAESAR SVM model gave good results on mutagenicity, with accuracy similar to the reproducibility of the experimental data. Strict quality check of the data has been done, and validation with about 850 compounds. The combination of CAESAR SVM model and ToxTree further reduced *false negative* rate.

We acknowledge Q. Chaudry, J. Cotterill (CSL, UK) and A. Chana (Mario Negri, Italy) for checking chemical structures.
We acknowledge financial contribution of the EC (project CAESAR).

