

Machine Learning Techniques for peptide optimization from sequence information

Mark R. Hansen, Jane Razumovskaya, Jason Hodges, Hugo O Villar
Altoris, Inc., San Diego, CA

www.althoris.com

info@althoris.com

Background

The increased interest in peptides, antibodies and other biopolymers as therapeutic agents has made obvious that a significant gap exists in research informatics for this class of drugs. Most of the tools available for peptide design aim to deploy molecular modeling techniques that can be powerful but as datasets become larger and more complex, the tools available to visualize, manage and organize the data are deficient. The problem is particularly acute for peptides where the common use of unnatural amino acids or chemical modifications precludes some of the techniques that are used for protein therapeutics. There is a clear need for tools that facilitate the discovery of relations between sequence and data. Two general types of methods are needed. First, techniques for **exploratory data analysis**, where we aim to discover relations in available data. Second **predictive analytics** techniques that aim to develop models with predictive power. Exploratory data analysis precedes predictive analytics, since a clear understanding of the available data is needed prior to embarking in the development of predictive models.

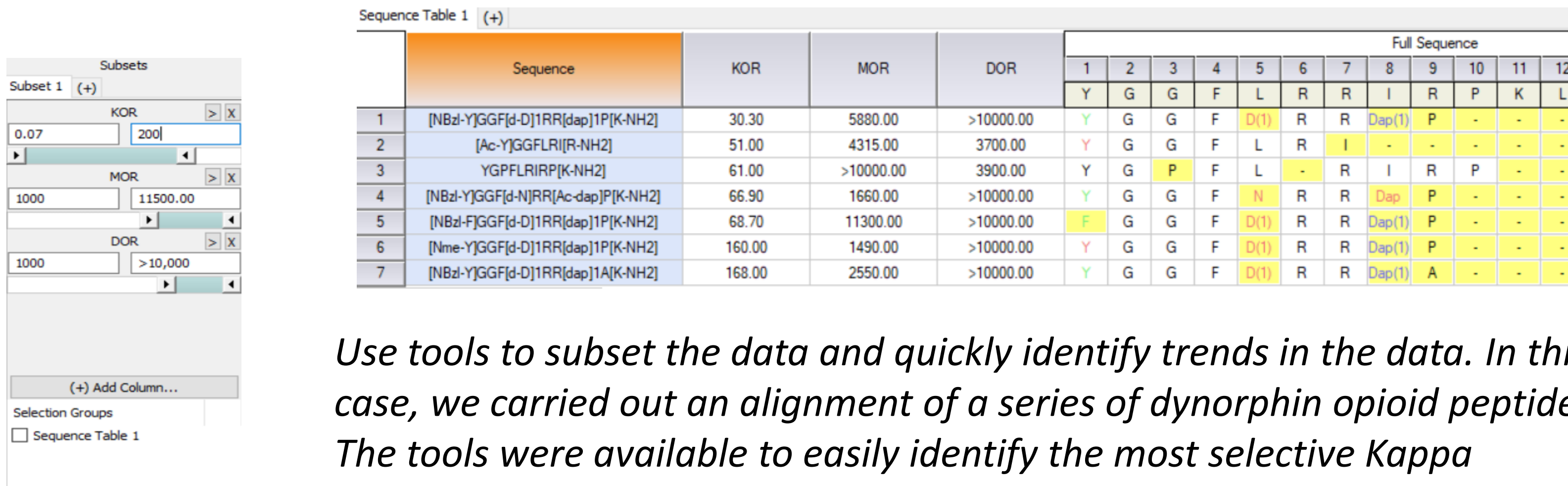
Exploratory Data Analysis for Biopolymers

In the past we have developed a tool to discover relations between data and sequences, a sequence activity relationship tool, named **SARvision|Biologics**.

Accurate Sequence Alignments are a critical aspect to the analysis and the predictive capabilities. *If the starting alignments are incorrect, then all subsequent studies will be inaccurate.*

Pairwise template based or Multisequence Alignment
Substitution matrices (PAM, Blosom, Identity, **Custom**)
Manual adjustments

A spreadsheet that relates aligned sequences to data provides the first layer of data analysis:

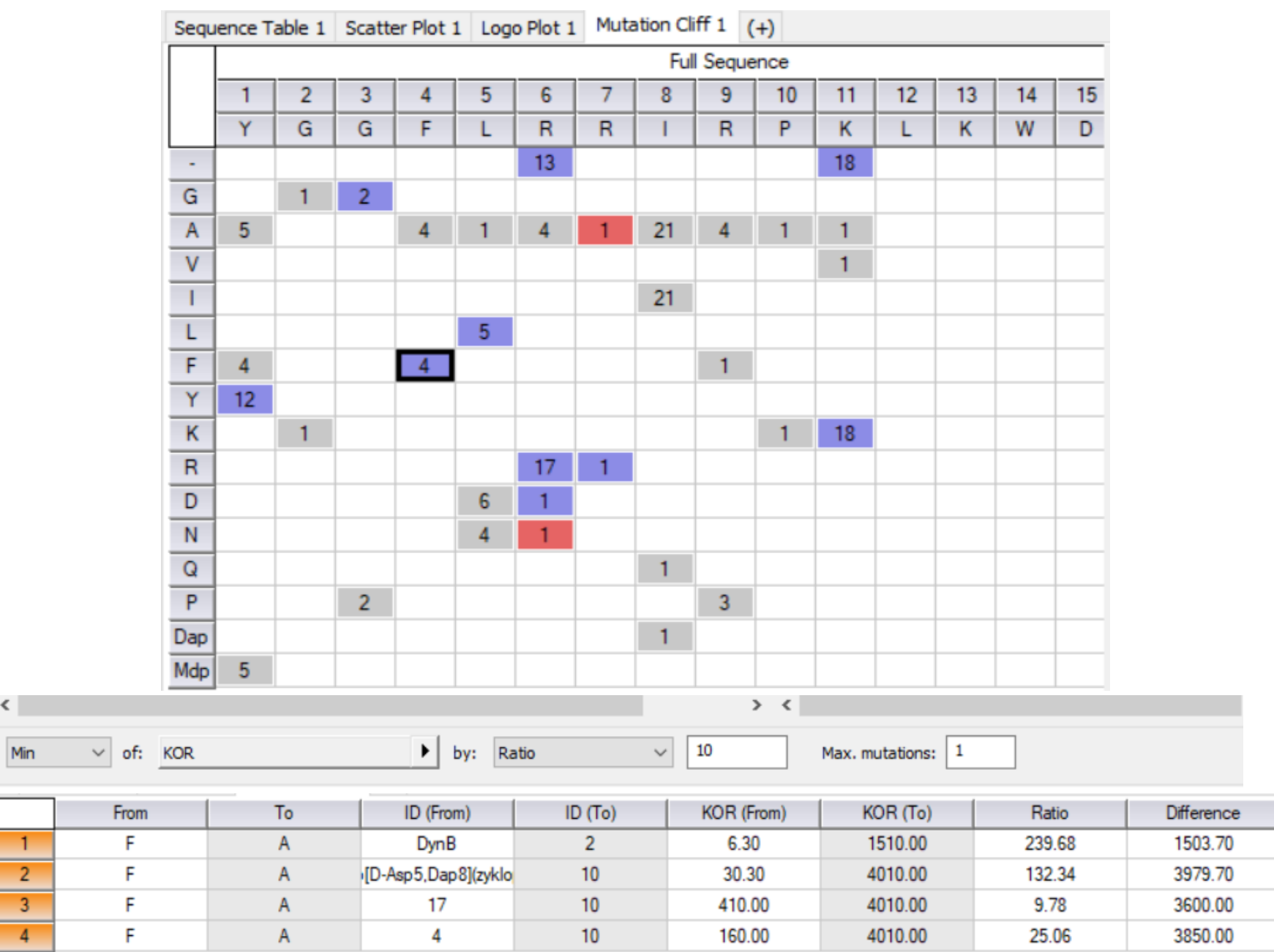
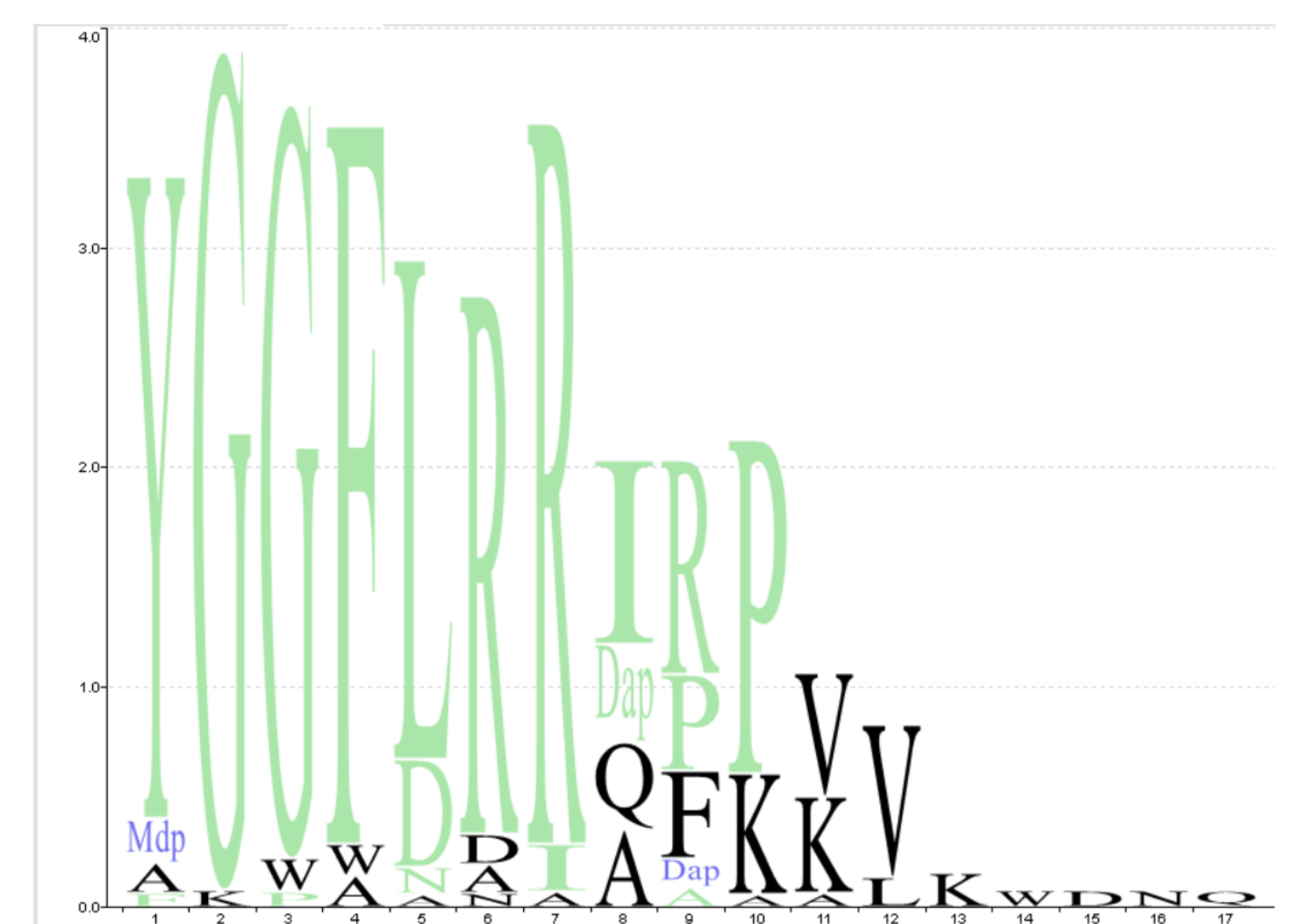


Use tools to subset the data and quickly identify trends in the data. In this case, we carried out an alignment of a series of dynorphin opioid peptides. The tools were available to easily identify the most selective Kappa peptides.

Tools to Identify Critical Residues are important to develop a, understanding of residues that are critical for activity. Models used for predictive techniques should reflect the observations made in the exploratory analysis of data.

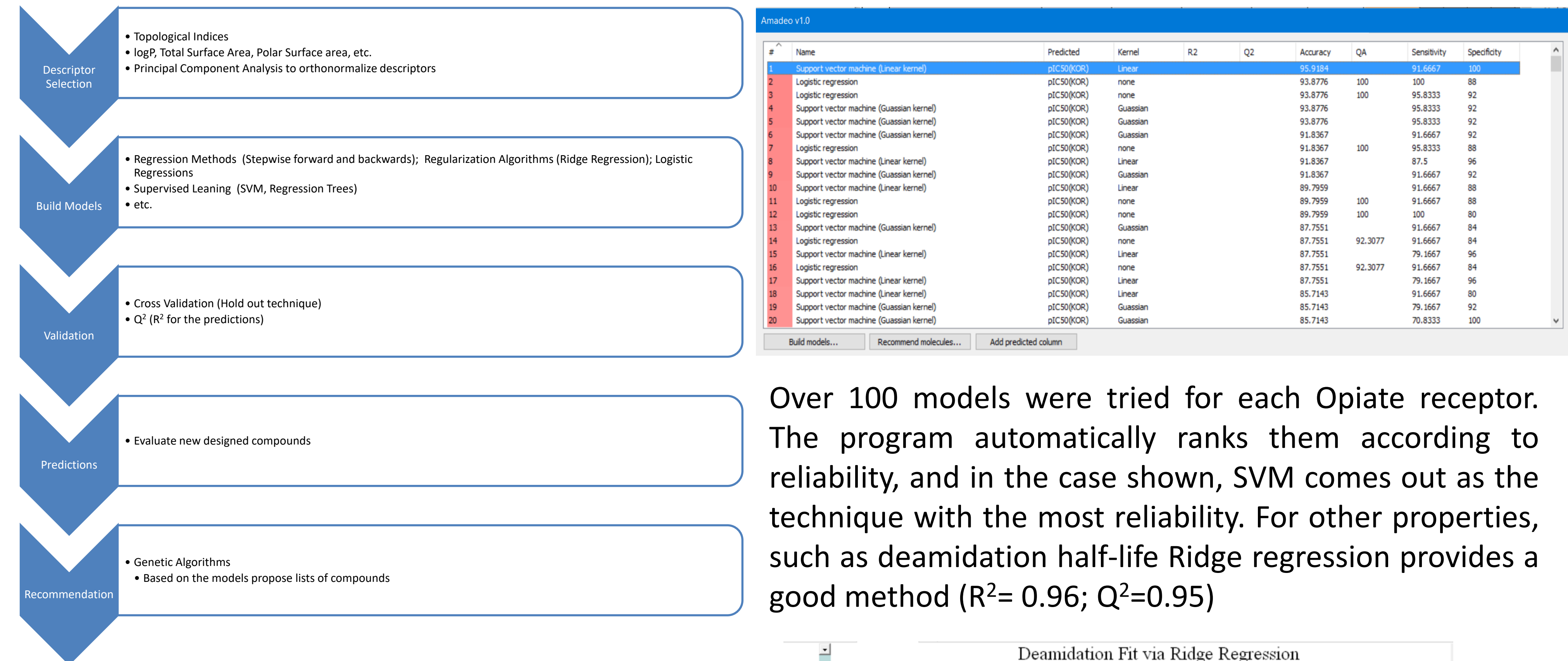


Different tools have to be used simultaneously to identify trends in the SAR data. Simple scatter plots, LogoPlots and Mutation Cliffs are some of the tools that can be used to that end. The LogoPlot shows in green the frequency of residues at each position of the KOR selective ligands. In position 5 there is some variability, while positions 2 and 4 have no variability, for the set selected. Mutation cliffs show where mutations in a position happen to have a significant effect in activity. For example in 4 compounds F->A mutations resulted in more than a 10x change in KOR affinity.

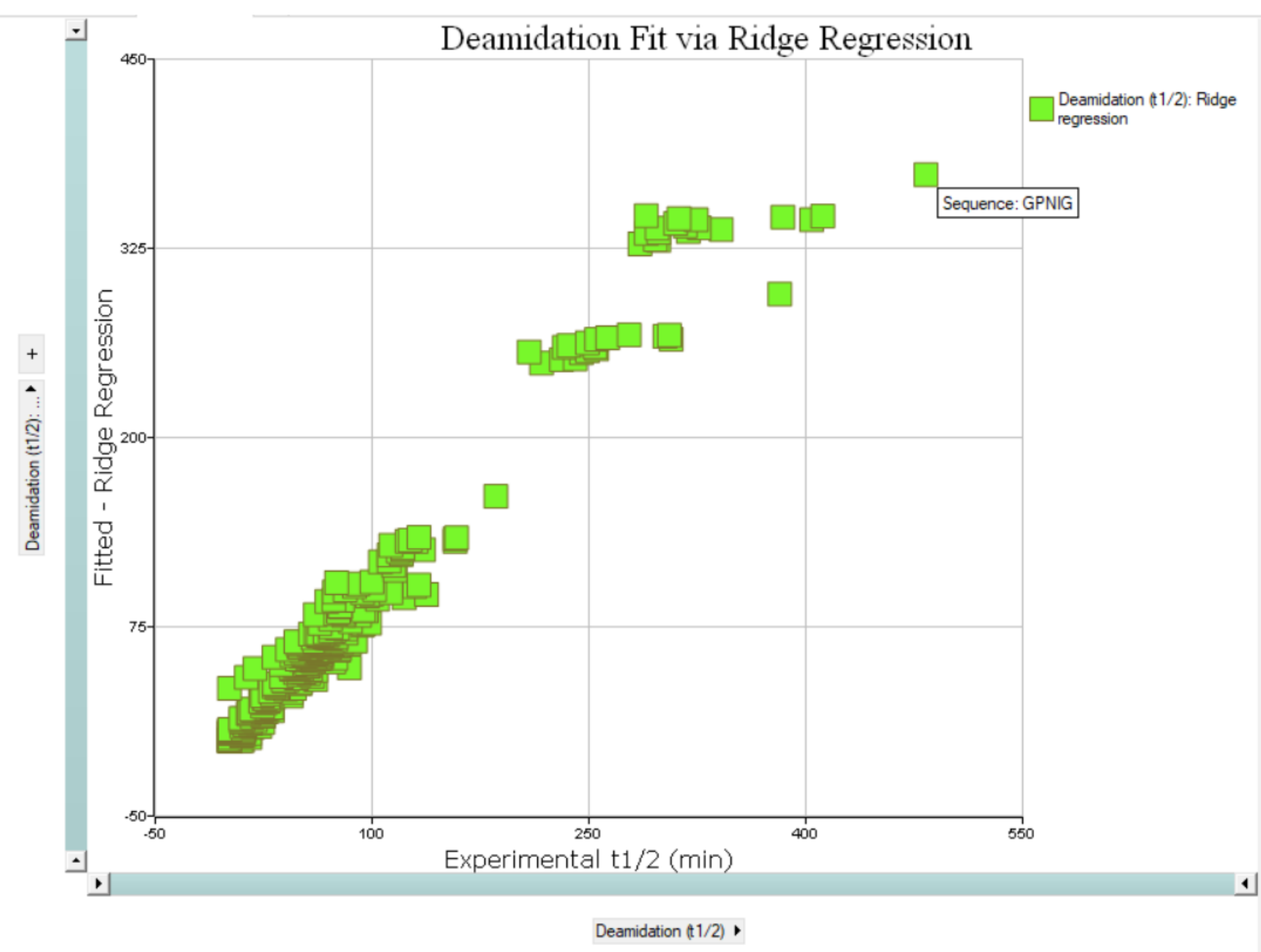


Machine Learning for Peptides

Machine Learning methods can be powerful when applied to the optimization of a compound. Unfortunately, the most appropriate methods for each property to be optimized can be different and require some expertise in their use. A computer system that automatically examines the utility of different machine learning techniques for a dataset and selects the most predictive methods for different properties of clinical interest for peptides was created. The program automatically will **LEARN** about your data, **BUILD** and **SELECT** the most predictive models, and make **RECOMMENDATIONS** as to what compounds to make next.



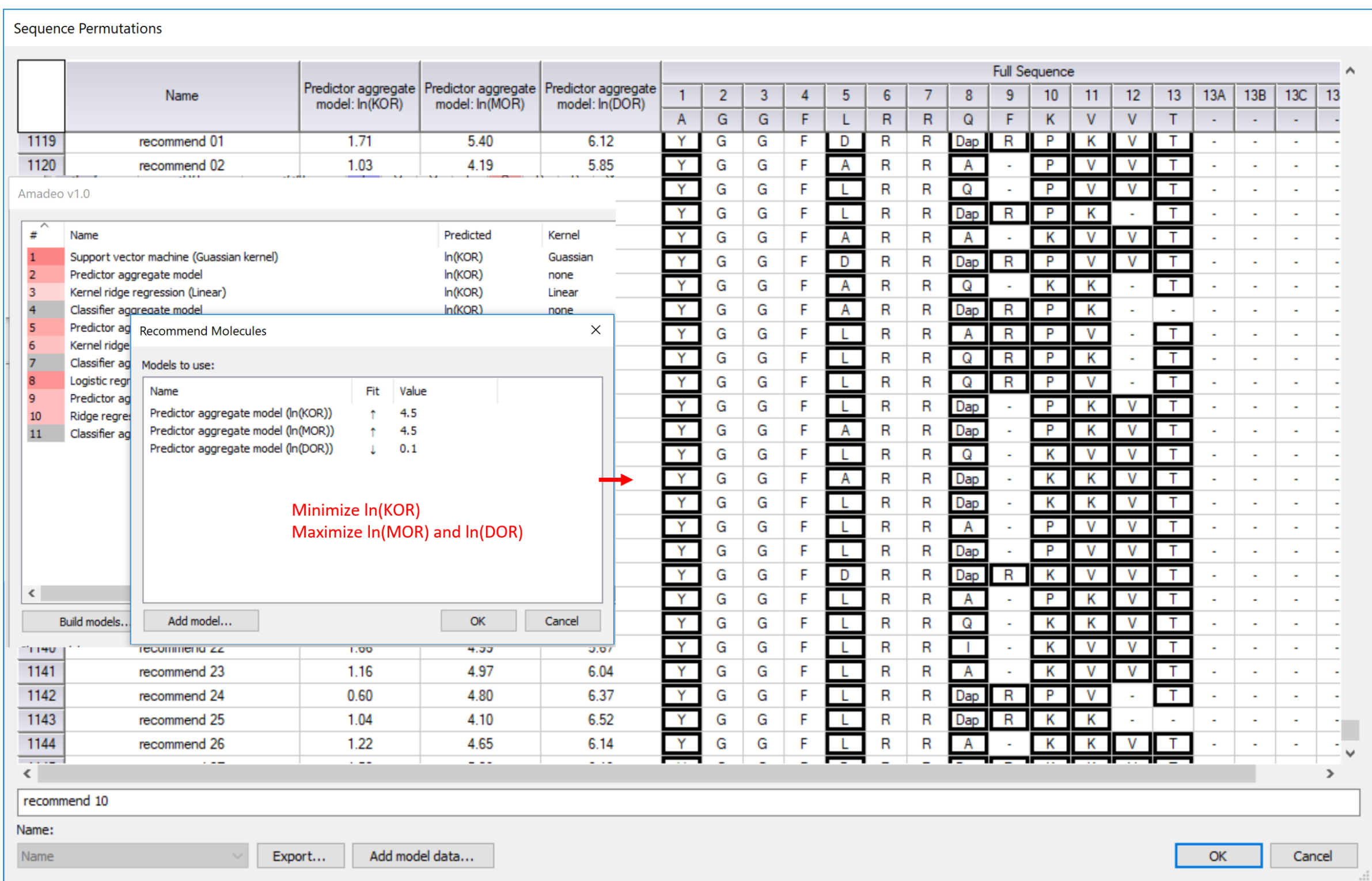
Over 100 models were tried for each Opiate receptor. The program automatically ranks them according to reliability, and in the case shown, SVM comes out as the technique with the most reliability. For other properties, such as deamidation half-life Ridge regression provides a good method (R²= 0.96; Q²=0.95)



The exploration of multiple methods and monomer descriptors is now fully automated, which makes these methods available to the bench scientist who may not be familiar with machine learning and artificial intelligence algorithms.

In some cases, the most predictive models are those that result from combining the responses of several models. The program we describe aims to optimize both. The program categorizes the predictive power of the different techniques so that it can recommend molecules that satisfy several parameters simultaneously.

Note that the program is able to handle numerical and categorical data, such as is the case for aggregation.



Future Directions

The combination of these best in class algorithms provides an avenue to solve the multifactorial problem of peptide optimization. Our ultimate goal is to develop a decision support system that guides the optimization of peptides towards the definition of a clinical candidate minimizing the number of peptides that need to be evaluated and useful to non-experts. Still some open questions remain:

- What is the best way to combine the different models to ensure they are predictive?
- Are the properties we selected optimal ? How to incorporate 3D information?
- Can other sources of data be incorporated as filters that enhance the accuracy of the predictions?

We are currently seeking projects that may benefit from the use of these tools, and that enable us to improve on them.

References

1. SARvision|Biologics, Altoris Inc., San Diego, CA www.chemapps.com
2. Hansen, Villar & Feyfant, J. Chem. Inf. Model. 2013, 53, 2774–2779
3. Joshi, Murray & Aldrich Biopolymers. 2017 Sep;108(5) ; Ramos-Colon et al. J Med Chem. 2016 59:10291; Joshi, Murray & Aldrich J Med Chem. 2015 58:8783; Patkar, Murray & Aldrich J Med Chem. 2009 52:6814