

# A predictive alignment-free model based on a new logistic regression-based method for feature selection in complete and partial sequences of Senecavirus A

Tatiana Flávia Pinheiro de Oliveira<sup>1</sup>, Marcos Augusto dos Santos<sup>2</sup>, Marcelo Fernandes Camargos<sup>3</sup>, Antônio Augusto Fonseca Júnior<sup>3</sup>, Aristóteles Góes-Neto<sup>4</sup>, Edel Figueiredo Barbosa Stancioli<sup>4</sup>,

*1 Departamento de Microbiologia, Instituto de Ciências Biológicas, UFMG, Ministério da Agricultura, Pecuária e Abastecimento*

*2 Departamento da Ciência da Computação, Instituto de Ciências Exatas, UFMG*

*3 Ministério da Agricultura, Pecuária e Abastecimento*

*4 Departamento de Microbiologia, Instituto de Ciências Biológicas, UFMG*

## Abstract

In 2015, there was an outbreak involving pig farms in six Brazilian states, whose single agent found and described for the first time in the country was the Senecavirus A (SVA), a virus belonging to the genus Senecavirus (Picornavirales, Picornaviridae). This viral family also houses the genus Aphthovirus, whose species type is the Foot-and-mouth disease virus (FMDV), agent of Foot-and-mouth disease, a highly infectious disease notifiable under the strict control of the Ministry of Agriculture, Livestock and Supply (MAPA) and the World Organisation for Animal Health (OIE). In the past few years, there has been a growing interest in the application of methods of linear algebra and statistics in data mining, social networks, machine learning, bioinformatics and information retrieval. Among these methods, logistic regression approach draws some special interest as it is a standard method for data classification using genome data and is the most frequently used method for disease prediction. We introduce a model that represents sequences as 6-nucleotide frequency vectors in R4096 and 3- amino acids frequency vectors in R800 and uses information of SVA and FMDV from the complete genome or amino acid sequences of the polyprotein of these viruses. In addition, partial sequences of nucleotides / amino acids of structural proteins (VP1, VP2, VP3 and VP4) of these viruses were used to build a new logistic regression-based method for classification. This new model allowed the assignment of values to parameters  $a_i^*$  that are associated with the frequency of a certain hexanucleotide or triplets codons of amino acids. Scrutinizing these parameters  $a_i^*$  unveiled that the most positive value may be related to important target sites of key virus proteins. Thus, this methodology was able to predict key regions in Senecavirus A, which can be important in studies of viral replication mechanism or in the development of diagnostic kits.

Funding: LANAGRO-MG ; FAPEMIG; CNPq