

QSAR, molecular docking Studies of 6-(Amino methyl)-5-(2,4-dichlorophenyl)-7-methylimidazo[1,2-a] pyrimidine-2- carboxamides as Potent, Selective Dipeptidyl Peptidase-4 (DPP4) Inhibitors

Abstract:

Type 2 diabetes (T2DM) is a metabolic disorder disease and DPP-4 inhibitors are a class of oral hypoglycemics that block the enzyme dipeptidyl peptidase-4 (DPP-4). DPP-4 inhibitors reduce glucagon and blood glucose levels and don't have side effects such as hypoglycemia or weight gain. In this paper a series of imidazolopyrimidine amides analogues as DPP4 inhibitors were applied for quantitative structure–activity relationship (QSAR) analysis. A collection of chemometric methods such as multiple linear regression (MLR), factor analysis-based multiple linear regression (FA-MLR), principal component regression (PCR) , genetic algorithm for variable selection-MLR (GA-MLR) and partial least squared combined with genetic algorithm for variable selection (GA-PLS), were conducted to make relations between structural features and DPP4 inhibitory of a variety of imidazolopyrimidine amides derivatives. GA-PLS represented superior results with a high statistical quality ($R^2 = 0.94$ and $Q^2 = 0.80$) for predicting the activity of the compounds. Docking studies of these compounds reveals and confirms that, compounds 15, 18, 25, 26, and 28 are introduced as good candidates for DPP-4 inhibitors were introduced as a good candidate for DPP-4 inhibitory compounds.

Keywords:

Imidazo pyrimidine derivatives, DPP-4 inhibitors, QSAR, Molecular docking

29 **1.Introduction**

30 Diabetes Mellitus (DM) is a metabolic disorder disease that body doesn't have the ability to
31 produce insulin or is resistant to insulin so it cannot function properly. Dipeptide peptidase 4
32 inhibitors are a new therapy Target that does not complicate previous medications such as
33 hypoglycemia, weight gain and cardiovascular risk. DPP-4 is a membrane protease that has a
34 specific selectivity on the secretion of incretins hormones in fact, these drugs can effective in
35 controlling the secretion of insulin and reducing glucagon secretion on hemostasis of glucose. So
36 reducing glucagon secretion and can also affect the process of gluconeogenesis in the liver.
37 Therefore, by inhibiting this physiological pathway in the body, effectively reduce blood glucose
38 levels.

39 The quantitative structure-activity relationship (QSAR) research field provides medicinal
40 chemists with the ability to predict drug activity by mathematical equations which construct a
41 relationship between biological activity of the molecules and descriptors [1, 2]. These
42 mathematical equations are in the form of $y = Xb + e$ that describe a set of predictor variables (X)
43 with a predicted variable (y) by means of a regression vector (b) [3]. The most important step in
44 building QSAR models is the appropriate representation of the structural and physicochemical
45 features of structures [4-10]. These features called molecular descriptors are the ones with
46 higher impact on the biological activity of interest. Nowadays, a wide range of descriptors are
47 being used in QSAR studies which can be classified into different categories according to the
48 Karelson approach including; constitutional, geometrical, topological, quantum, chemical and so
49 on [8].Hyperchem and Dragon are two well-known computational softwares provide us more
50 than 1000 of these descriptors [11-12].There are different variable selection methods available
51 including; stepwise multiple linear regression (MLR), genetic algorithm (GA), principal
52 component or factor analysis (PCA) and so on. The mathematical relationships between
53 molecular descriptors and activity are used to find the parameters affecting the biological activity
54 and/or estimate the property of other molecules.

55 Here, we consider the DPP4 inhibitory activity of a novel series of imidazolopyrimidine
56 amides which have been recently designed and synthesized by W. Meng [13]. Our research show
57 that these series of compounds don't evaluate for QSAR studies. Different statistical methods
58 were applied to model the relationship between the structural features and the DPP-4 inhibitory
59 activity of the studied compounds. These methods are: (i) multiple linear regression (MLR), (ii)

60 principal component regression (PCR), MLR with factor analysis as the data pre-processing step
61 for variable selection (FA-MLR) (iii), genetic algorithm-multiple linear regression (GA-MLR)
62 (iv), genetic algorithm - partial least squares (GA-PLS) (v). Molecular docking simulation
63 technique was also performed on twenty-nine compounds to reach the details molecular binding
64 models for these compounds interacting with the key active site DPP-4 inhibitors.

65

66 **2. Materials and methods**

67 **2.1. Data set**

68 The biological activity was used in this study, were the DPP-4 inhibitory activity of a set of
69 thirty-one imidazolepyrimidine amides derivatives (13), which were designed, synthesized and
70 evaluated for their ability as potential treatments for type II diabetes. The structural features and
71 biological activities of these compounds are listed in Table 1. The biological data were converted
72 to logarithmic scale (pIC50) and then used for subsequent QSAR analysis as dependent variable.

73

74 [Table 1. near here]

75

76 **2.2 Molecular descriptors**

77 All structures were generated with HyperChem program (Hyper-cube Inc., Version 8.0.3) [11]
78 and optimized by MM+ method and then semi-empirical AM1 method in hyperchem software.
79 The molecular structures were optimized using the Polak-Ribiere algorithm until the root mean
80 square gradient was $0.01 \text{ kcal mol}^{-1}$. Some chemical parameters including molar volume (V),
81 molecular surface area (SA), hydrophobicity (logP), hydration energy (HE) and molecular
82 polarizability were calculated by using Hyperchem software. The resulted geometry was
83 transferred into Dragon program package, which was developed by Milano Chemometrics and
84 QSAR Group. Dragon software (version 5.5) [12] calculated different topological, geometrical,
85 charge, empirical and constitutional descriptors for each molecule. 2D autocorrelations
86 aromaticity indices, atom-centered fragments and functional groups were also calculated by
87 dragon software.

88 In the case of docking procedure, each optimized structures in HyperChem 8.0.3 program were
89 there after converted to PDBQT using MGLtools 1.5.6 [14]. The three dimensional crystal
90 structure of dipeptidyl peptidase iv human (PDB ID:5j3j) were retrieved from protein data bank

91 [15]. Co-crystal ligand molecules were excluded from the structures and the PDBs were
92 corrected in terms of missing atom types by modeller9.12 [16]. An in house application
93 (MODELFACE) was used for generation of python script and running modeler software [17].
94 Subsequently, the enzymes were converted to PDBQT and gasteiger partial charges were added
95 using MGLtools1.5.6.

96

97 **2.3. Data screening and model building**

98 The calculated descriptors were collected in a data matrix, D whose number of rows and columns
99 were the number of molecules and descriptors, respectively. First the descriptors were checked
100 for constant or near constant values and those detected were removed from the original data
101 matrix. The correlated descriptors with each other's and with the activity data were determined
102 and removed from pool of descriptors.

103 Five different methods were used: (1) stepwise-multiple linear regression (2) MLR with factor
104 analysis as the data pre-processing step for variable selection (FA-MLR), (3) principal
105 component regression analysis and (4) genetic algorithm- multiple linear regression (GA-MLR)
106 (5) genetic algorithm- partial least squares (GA-PLS).

107 MLR with stepwise selection and elimination of variables was applied for developing QSAR
108 models by using SPSS software (SPSS Inc., version 21). The resulted models were validated by
109 leave-one out cross-validation procedure by using MATLAB software version 2014. However,
110 this procedure did not produce good results and therefore we used genetic algorithm (GA-PLS)
111 to select the best variables. FA-MLR was performed on the dataset. Factor analysis was used to
112 reduce the number of variables. Principal component regression analysis was also tried for the
113 dataset along with FA-MLR. With PCRA collinearities among X variables are not a disturbing
114 factor and the number of variables included in the analysis may exceed the number of
115 observations [18]. In this method, factor scores, as obtained from FA, are used as the predictor
116 variables [19]. In PCRA, all descriptors are assumed to be important while the aim of factor
117 analysis is to identify relevant descriptors. Partial least squares (PLS) linear regression is a recent
118 technique that generalizes and combines features from principal component analysis and
119 multiple regressions. PLS is a method suitable for overcoming the problems in MLR related to
120 multicollinear or over-abundant descriptors [20]. Application of PLS method thus allows the
121 construction of larger QSAR equations while still avoiding over-fitting and eliminating most

122 variables. This method is normally used in combination with cross-validation to obtain the
123 optimum number of components [21]. The PLS regression method used was the NIPALS-based
124 algorithm existed in the chemometrics toolbox of MATLAB software (version 8.0.3.532 Math
125 Work Inc.).

126

127 **2.4. Docking procedures**

128 An in house batch script (DOCK-FACE) for automatic running of AutoDock 4.2 was used to
129 carry out the docking simulations [22] in a parallel mode [23]. To prepare the receptor structure,
130 the three dimensional crystal structure of Dipeptidyl Peptidase-4 (PDB ID: 5j3j) was acquired
131 from Protein Data Bank (PDB data base; <http://www.rcsb.org>) [24] and water molecules and co-
132 crystal ligand were removed from the structure. The PDB were then checked for missing atom
133 types with the python script as implemented in MODELLER 9.17 [25]. The ligand structures
134 were made by Hyper Chem software package (Version 7, Hypercube Inc). For geometry
135 optimization, Molecular Mechanic (MM⁺), followed by semi empirical AM1 method was
136 performed. The prepared Ligands were given to 100 independent genetic algorithm (GA) runs.
137 150 population size, a maximum number of 2,500,000 energy evaluations and 27,000 maximum
138 generations were used for Lamarckian GA method. The grid points of 30, 30, and 30 in x-, y-,
139 and z directions 20.3, 3.7 and 51.3 were used. Number of points in x, y and z was and 65
140 respectively. All visualization of protein ligand interaction was evaluated using VMD software
141 [26]. Cluster analysis was performed on the docked results using a root mean square deviation
142 (RMSD) tolerance of 2.4 Å.

143

144

145

146

147

148 3. Results and Discussion

149 The structural feature and the experimental DPP-4 inhibitory activity (represented as pIC_{50}) of
150 the molecules used in this study are shown in Table 1. To obtain the effects of the structural
151 parameters of the investigated derivatives on their DPP-4 activity, QSAR analysis was
152 performed with various molecular descriptors. Among the different chemometric tools available
153 for modeling the relationship between the biological activity and molecular descriptors, five
154 methods (i.e., stepwise MLR, PCR, FA-MLR, GA-MLR and GA-PLS) were applied and
155 compared here. The calculated descriptors from whole molecular structures are briefly described
156 in Table 2.

157
158 [Table 2. near here]
159

160 3.1. MLR models for subset of molecules

161 Firstly, separate stepwise selection-based MLR analyses were performed using different types of
162 descriptors, and then, a MLR equation was obtained utilizing the pool of all calculated
163 descriptors. First principal component analysis was done to detect outlier data and was drawn
164 PC1 on PC2 (Figure 1), as it can show the molecule number of 1 and 16 are outlier data so
165 omitted. Then Kennard stone algorithm was used to divide data set to calibration and prediction
166 set. MLR models with maximum number of variables of 5 were selected. Statistical parameters
167 such as correlation coefficient (R^2), correlation coefficient for test set ($R^2_{\text{test set}}$ or R^2_{predic}),
168 standard error of regression (SE), and Fisher ratio (F) at specified degrees of freedom, leave-one-
169 out cross-validation correlation coefficient (Q^2) was shown in Table 3. Equation 1 was selected
170 as the best equation in the MLR model because of its greatest statistical parameters. The selected
171 variables demonstrate that 2D-autocorrelation (MATS1m), constitutional (Ms), topological
172 charge indices (GGI5), topological (DELS), 3D-MORSE descriptors (Mor25m) effect on the
173 inhibitory activity of the studied compounds.

174 A small difference between the conventional and cross-validate correlation coefficients of the
175 different MLR equations (Table 4) reveals that none of the models are over fitted, which can be
176 partially attributed to the absence of collinearity between the variables in one hand and use of no
177 extra variables on the other hand. Equation 1 (as the best equation in this series) could explain
178 91% of the variance and predict 84% of the variance in $(-\log IC_{50})$ data. All of the descriptors

179 that used in this equation have positive effect on DPP-4 inhibitory expect DELS as topological
180 descriptors. Figure 2 shows the plots of linear regression predicted versus experimental value of
181 the DPP4 inhibitory activity of ligand. The plots for this model show to be more convenient with
182 $R^2_{cv}= 0.84$.

183 [Table 3. near here], [Table 4. near here]

184 [Figure 1. near here], [Figure 2. near here]

185

186 3.2. PCR Analysis

187 When factor scores were used as the predictor parameters in a multiple regression equation
188 (Table 5), a predictive QSAR model with factor scores of 1, 2, 3 and 7 as input variables, was
189 obtained (Table 3, Equation 2). This equation shows statistical quantities similar to those
190 obtained by the FA-MLR method.

191 Considering this information in modeling, it may apparently increase the model variances (i.e.,
192 R^2) but they are useful for prediction. Figure 2 shows the plots of linear regression predicted
193 versus experimental value of the DPP-4 inhibitory activity of ligand. The plots for this model
194 show to be more convenient with $R^2_{cv}= 0.75$.

195 [Table 5. near here]

196 3.2. FA-MLR analysis

197 FA-MLR was performed on the dataset. Factor analysis (FA) was used to reduce the number of
198 variables and to detect structure in the relationships between them. This data-processing step is
199 applied to identify the important predictor variables and to avoid collinearities among them.
200 Principle component regression analysis, PCRA, was tried for the dataset along with FA-MLR.
201 With PCRA collinearities among \mathbf{X} variables are not a disturbing factor and the number of
202 variables included in the analysis may exceed the number of observations [27]. In this method,
203 factor scores, as obtained from FA, are used as the predictor variables [28]. In PCRA, all
204 descriptors are assumed to be important while the aim of factor analysis is to identify relevant
205 descriptors. Table 5 shows the two factor loadings of the variables (after VARIMAX rotation)
206 for the compounds tested against dipeptidyl peptidase 4 inhibitors'. As it is observed, about 77%
207 of variances in DPP4 inhibitors' could be explained by the selected two factors. It is observed;
208 about 0.67 of variances in the original data matrix can be explained by selected 2 factors,
209 MATS1m as 2D-autocorrelation descriptors and Ms as Constitutional descriptors. And also have

210 weak predicted variance in DPP4 inhibitory. Figure 2 shows the plots of linear regression
211 predicted versus experimental value of the DPP4 inhibitory activity of ligand. The plots for this
212 model show to be more convenient with $R^2_{cv} = 0.53$.

213

214 **3.3.GA-MLR analysis**

215 Genetic algorithm technique was employed as a selection tool to select the most relevant
216 descriptors with respect to an objective function. The genetic algorithm (GA) starts with the
217 creation of a population of randomly generated parameter sets. the parameters set used for the
218 GA includes population size (160), initial terms 18%, max generation (250) and %convergences
219 (90%), These selected subsets of variables are further evaluated by their fitness to predict
220 inhibitory activity values. multiple linear regression analysis was performed on the training set
221 and then, evaluated by test set. Using genetic algorithm-multiple linear regression (GA-MLR)
222 analysis resulted in the development of a predictive QSAR model with four descriptors with the
223 following equation:

$$224 \text{PIC50} = 2.072\text{GGI7} ((\pm 0.676)) + 4.427\text{Ms}((\pm 0.678)) + 8.047\text{BELm6}(\pm 1.753) - 0.453\text{MOr27u}(\pm 0.187) - \\ 225 14.411(\pm 3.275)$$

226 The statistical parameters of GA-MLR model are shown in Table 3. and could explain 94% of the
227 variance and predict 88% of the variance in $(-\log IC_{50})$ data. This equation describes the effect of
228 GGI7 (Topological charge indices), Ms (Constitutional), BELm6 (Burden Eigenvalues) and
229 MOR27U (3-D Morse Descriptors) in dpp4 inhibitory. All the descriptors have positive
230 coefficient except MOR27u and indicated that increase this descriptor (MOR27u) could result in
231 decreasing PIC50. Figure 2 shows the plots of linear regression predicted versus experimental
232 value of the dpp4 inhibitory activity of ligand. The plots for this model show to be more
233 convenient with $R^2_{cv} = 0.88$.

234

235 **3.4.GA-PLS analysis**

236 In PLS analysis, the descriptors data matrix is decomposed to orthogonal matrices with an inner
237 relationship between the dependent and independent variables. Therefore, unlike MLR analysis,
238 the multi collinearity problem in the descriptors is omitted by PLS analysis. Because a minimal
239 number of latent variables are used for modeling in PLS; this modeling method coincides with
240 noisy data better than MLR. In order to find the more convenient set of descriptors in PLS

241 modeling, genetic algorithm was used. To do so, many different GA-PLS runs were conducted
242 using different initial set of populations.

243 The data set ($n = 29$) was divided into two group: calibration set ($n = 20$) and prediction set ($n =$
244 9). Given 20 calibration samples; the leave-one out cross-validation procedure was used to find
245 the optimum number of latent variables for each PLS model.

246

247 [Table 6. near here]

248

249 The most convenient GA-PLS model that resulted in the best fitness contained 9 indices, five of
250 them being those obtained by MLR. The PLS estimate of coefficients for these descriptors are
251 given in Figure 3. As it observed, a combination of Constitutional, Topological, Connectivity
252 indices, 2D-autocorrelation, Edge adjacency indices, Topological charge indices, 3-D Morse
253 Descriptors, WHIM Descriptors have been selected by GA-PLS to account the Dipeptidyl
254 Peptidase-4 (DPP4) inhibitory activity of imidazole derivatives. The resulted GA-PLS model
255 possessed a high statistical quality $R^2 = 0.94$ and $Q^2 = 0.80$. The predictive ability of the model
256 was measured by applying to 10 external tests set molecules. The squared correlation coefficient
257 for prediction was 0.95 and standard error of prediction was 0.49. The values of pIC_{50} using GA-
258 PLS model (refined from cross-validation or external prediction set) are shown in Table 1. This
259 figure 3 describes the effect of Ms (Constitutional), X0A (Connectivity indices), MATSIM(2D-
260 autocorrelation), EEig09d, EEig13d (Edge adjacency indices), GGi5(Topological charge
261 indices), Mor25m (3-D Morse Descriptors), E2M (WHIM Descriptors) and DELS (Topological)
262 on inhibitory DPP4 activity. And also describe that X0A, EEig13d and DELS have negative
263 coefficient on DPP4 inhibitory but the other descriptors have positive effect on DPP4 activity.
264 Figure 2 shows the plots of linear regression predicted versus experimental value of the DPP4
265 inhibitory activity of ligand. The plots for this model show to be more convenient with $R^2_{cv} =$
266 0.80.

267 In order to investigate the relative importance of the variable appeared in the final model
268 obtained by GA-PLS method, variable important in projection (VIP) was employed [29]. VIP
269 values reflect the importance of terms in PLS model. According to Erikson et al. X-variables
270 (predictor variables) could be classified according to their relevance in explaining y (predicted

271 variable), so that $VIP > 1.0$ and $VIP < 0.8$ mean highly or less influential, respectively, and $0.8 <$
272 $VIP < 1.0$ means moderately influential [30].

273

274

275 [Figure 4. Near here]

276

277 The VIP analysis of PLS equation is shown in Figure 4. VIP analysis shows that M_s which is
278 constitutional descriptors, $X0A$ as Connectivity indices descriptors, $MATS1m$ which is 2D-
279 autocorrelation and $GGI5$ which is topological charge indices parameter, are the most important
280 indices in the QSAR equation derived by PLS analysis. In addition, the other descriptors have
281 been found to be low influential parameters.

282

283 **3.5. Robustness and applicability domain of the models**

284 Leverage is one of standard methods for this purpose. Warning leverage (h^*) is another criterion
285 for interpretation of the results. The warning leverage is, generally, fixed at $3k/n$, where n is the
286 number of training compounds and k is the number of model parameters. A leverage greater than
287 warning leverage h^* means that the predicted response is the result of substantial extrapolation of
288 the model and therefore may not be reliable [31]. The calculated leverage values of the test set
289 samples for different models and the warning leverage, as the threshold value for accepted
290 prediction, are listed in Table 6. As seen, the leverages of all test samples are lower than h^* for
291 all models. This means that all predicted values are acceptable.

292 [Table.6 near here]

293

294 **3.6. Docking Study**

295 Docking is frequently used to predict the binding orientation of small molecule drug candidates
296 to their protein targets in order to predict the affinity and activity of the small molecule. Hence
297 docking plays a great role in the rational design of drugs. Here, docking studies were carried out
298 on our compounds to find their binding site, binding modes and the best direction on the base of
299 their binding energy. Having completed the docking process, the protein– ligand complex was
300 analyzed to investigate the type of interactions. The conformation with the lowest binding energy
301 was considered as the best docking result in each case resource.

302 As it was shown in Table 1, Compounds 15, 18, 25, 26 and 28 based on their highest docking
303 binding energy can be a good candidate for DPP-4 inhibitors.

304 On the other hand, promising results such as the ligand-receptor binding site and binding modes
305 were obtained from docking analyses. The results for each ligand were compared to its
306 corresponding co-crystal ligand.

307 [Figure 5. near hear]

308 On the other hand, promising results such as the ligand-receptor binding site and binding modes
309 were obtained from docking analysis. The results for each ligand were compared to its
310 corresponding co-crystal ligand. a hydrogen bond acceptor interaction between amino group of
311 co-crystal ligand (HL1) and Glu 166, Glu 167 and Tyr 623 of the receptor (salt bridge). Trifluoro
312 phenyl group was occupying S1 hydrophobic pocket of dpp4 inhibitors with val 627, His701,
313 Val617, Tyr 592, Tyr 627 and Tyr506 residue of receptor [32]. Dimethoxy phenyl group of co-
314 crystal ligand π - π interaction with Phe 318 of the receptor figure 5a. Hydrogen bindings between
315 docked potent agents such as 25 and the dpp4 receptor (5j3j) figure 5b. dichloro phenyl group π -
316 π interaction with Tyr 508 and were occupying S1 pocket. Hydrogen bonding interaction
317 between methyl of imidazole pyrimidine of 25 molecules and Trp 590 of receptor.

318

319 4. Conclusion

320 In this study, five different QSAR modeling methods, MLR, FA-MLR, PCR, GA-PLS and GA-
321 MLR were used in the construction of a QSAR model for DPP4 inhibitory of
322 imidazolopyrimidine amides and the resulting models were compared. The reliability, accuracy
323 and predictability of the proposed models were evaluated by root mean square error of cross-
324 validation (RMSECV) and cross- validation, the root mean square error of prediction (RMSEP).
325 Results confirm that among the applied models, the GA-PLS is superior for the prediction of the
326 pIC_{50} of imidazolopyrimidine amides analogues. All models represent high goodness of fit
327 (measured by R^2), whereas that obtained from GA-PLS is significantly better than that of the
328 other models. The cross-validation statistics reported suggested that the higher prediction ability
329 of the GA-PLS model. This study suggests the importance of constitutional, topological,
330 connectivity indices, 2D-autocorrelation, edge adjacency indices, topological charge indices, 3D
331 Morse descriptors, WHIM Descriptors of molecules for imidazolopyrimidine amides derivatives.

332 Docking study reveals and confirms that, compounds 15, 18, 25, 26, and 28 are introduced as
333 good candidates for DPP-4 inhibitors.

334
335 Quantitative relationships between molecular structure and anti-cancer activity of isatin
336 derivatives were discovered by two chemometrics methods: MLR and GA-PLS. Different QSAR
337 models revealed that topological parameters (X3v and PJI2) have significant impacts on the anti-
338 cancer activity of the compounds. In this series a significant role of chemical (Polarizability and
339 HE), constitutional (Sv and nX) and geometrical parameters (G1 and SPAM) on the inhibitory
340 activity was observed. Using the pool of all types of calculated descriptors a new QSAR model
341 was derived for these compounds. In this model the importance of effects of topological (X3v
342 and PJI2) and functional groups parameter (nCs) on the cytotoxic activity was indicated. The
343 positive effects of the number of halogen atoms and the number of total secondary carbons, and
344 the negative effects of the number of secondary amides, and the number of ketones on the anti-
345 cancer activity was in agreement with previous SAR studies.. GA-PLS model showed the effects
346 of seven topological indices (X3v, PW4, PJI2, MSD, SEigZ, IC1 and BIC2), three constitutional
347 descriptors (Ss, Me and nBr) and one chemical parameter (Vol) on the cytotoxic activity of the
348 compounds. A comparison between the two statistical methods employed indicated that MLR
349 represented superior results. The resulted MLR model possessed a high statistical quality ($R^2 =$
350 0.92 and $Q^2 = 0.90$) for predicting the activity of the compounds.

351

352

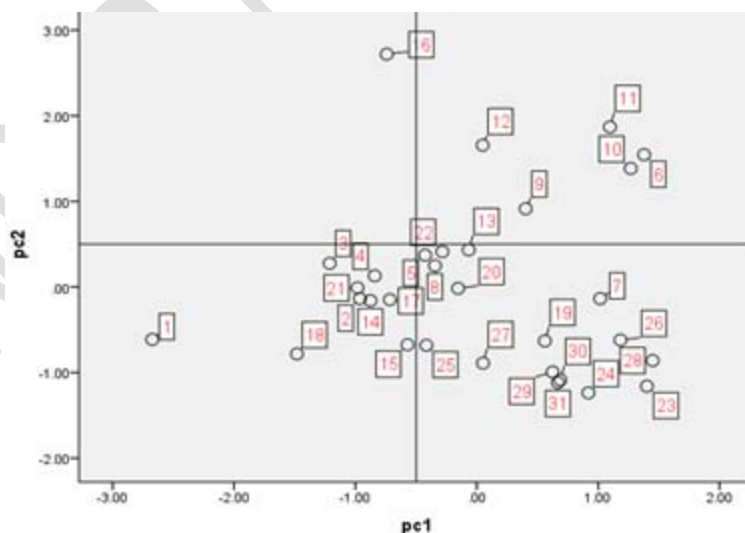
353 **References**

- 354 1. Hansch, C.; Hoekman D.; Gao, H. Comparative QSAR: Toward a Deeper Understanding of
355 Chemicobiological Interactions. *Chem. Rev.* **1996**, *96*, 1045-1076.
- 356 2. Hansch, C.; Maloney, P.P.; Fujita, T.; Muir, R.M. Correlation of Biological Activity of
357 Phenoxyacetic Acids with Hammett Substituent Constants and Partition Coefficients.
358 *Nature* **1962**, *194*, 178-180.
- 359 3. A. Fassihi, R. Sabet, QSAR Study of p56^{lck} Protein Tyrosine Kinase Inhibitory Activity of
360 Flavonoid Derivatives Using MLR and GA-PLS, *Int. J. Mol. Sci.* **9** (2008) 1876-1892.
- 361 4. Hansch, T. Fujita, ρ - σ - π Analysis. A method for the correlation of biological activity and
362 chemical structure, *J. Am. Chem. Soc.* **86** (1964) 1616-1626.
- 363 5. Sabet, R.; Fassihi, A.; Moeinifard, B., QSAR study of PETT Derivatives as Potent HIV-
364 Reverse Transcriptase Inhibitors. *J. Mol. Graph & Model.* **2009**; *28*: 146.
- 365 6. C. Hansch, D. Hoekman, H. Gao, Comparative QSAR: Toward a deeper understanding of
366 chemicobiological interactions, *Chem. Rev.* **96** (1996) 1045-1075.
- 367 7. R. Todeschini, V. Consonni, *Handbook of Molecular Descriptors*. Wiley-VCH, Weinheim,
368 2000.
- 369 8. Sabet, R.; Fassihi, A., QSAR study of Isatin analogues as in vitro anti-cancer agents. *Eur. J.*
370 *Med. Chem.* **2010**; *45*: 1113.
- 371 9. Sabet R.; Fassihi A.; Hemmateenejad B.; Saghaie L.; Miri R.; Gholami M.; Computer-aided
372 drug design of novel antibacterial 3-hydroxypyridine-4-ones: application of QSAR methods
373 based on the MOLMAP approach. *Journal of Computer-Aided Molecular Design.* **2012**;
374 26,349.
- 375 10. Karbakhsh, R.; Sabet, R.; Application of different chemometrics tools in QSAR Study of
376 Azolo-adamantanes against influenza A virus. **2011**; *6*,23.
- 377 11. Visit the Hyperchem official website at: <http://www.hyper.com>.
- 378 12. Todeschini, R. Milano Chemometrics and QSPR Group. <http://micchem.disat.unimib.it/>
- 379 13. Wei Meng, Robert P. Brigance, Hannguang J. Chao, Aberra Fura, Discovery of 6-(
380 Aminome thyl) -5-(2,4-dichlorophe nyl) -7-methy limidazo[1,2-a] pyrimidine-2-carbox
381 amides as Potent, Selective Dipeptidyl Peptidase-4 (DPP4) Inhibitor s. *J. Med. Chem.* **2010**,
382 *53*, 5620–5628.

- 383 14. Morris GM, Huey R, Olson AJ. Using AutoDock for Ligand-Receptor Docking. *Curr*
384 *Protoc Bioinformatics*. 2008; Chapter 8: Unit 8.14
- 385 15. Hikiş P, Szczupak Ł, Koceva-Chyła A, Oehninger L, Ott I, Therrien B, *et al.* Anticancer
386 and Antibacterial Activity Studies of Gold (I)-Alkynyl Chromones. *Molecules*. 2015;
387 20:19699-718
- 388 16. Eswar N, Webb B, Marti-Renom MA, Madhusudhan MS, Eramian D, Shen M-y, *et al.*
389 Comparative Protein Structure Modeling Using Modeller. *Curr Protoc Bioinformatics*.
390 2006; Chapter 5: Unit 5.6
- 391 17. Sakhteman A. PreAuposSOM, [https:// www.biomedicale.univ-paris5.fr/aupossom/](https://www.biomedicale.univ-paris5.fr/aupossom/)
- 392 18. Leardi, R. Application of Genetic Algorithm-PLS for Feature Selection in Spectral Data
393 Sets. *J. Chemometr.* **2000**, *14*, 643-655
- 394 19. Siedlecki, W.; Sklansky, J. On Automatic Feature Selection. *Int. J. Pattern Recog. Artif.*
395 *Intell.*, **1988**, *2*, 197-220.
- 396 20. Schmidti, H. Multivariate Prediction for QSAR. *Chemom. Intell. Lab. Sys.* **1997**, *37*, 125-
397 134
- 398 21. Hansch, C.; Kurup, A.; Garg, R.; Gao, H. Chem-bioinformatics and QSAR. A Review of
399 QSAR Lacking Positive Hydrophobic Terms. *Chem. Rev.* **2001**, *101*, 619-672.
- 400 22. Humphrey W, Dalke A, Schulten K. VMD: visual molecular dynamics. *Journal of*
401 *molecular graphics*. 1996;14(1):33-8.
- 402 23. Fereidoonnehad M, Faghieh Z, Mojaddami A, Sakhteman A, Rezaei Z. A Comparative
403 Docking Studies of Dichloroacetate Analogues on Four Isozymes of Pyruvate
404 Dehydrogenase Kinase in Humans. *Indian J Pharm Educ.* 2016;50(2):S32-S8.
- 405 24. Mirjalili BF, Zamani L, Zomorodian K, Khabnadideh S, Haghhighijoo Z, Malakotikhah Z, *et*
406 *al.* Synthesis, antifungal activity and docking study of 2-amino-4H-benzochromene-3-
407 carbonitrile derivatives. *Journal of Molecular Structure*. 2016; 1116:102-8.
- 408 25. Li Z, Gu J, Zhuang H, Kang L, Zhao X, Guo Q. Adaptive molecular docking method based
409 on information entropy genetic algorithm. *Applied Soft Computing*. 2015; 26:299-302.
- 410 26. Feng J, Ablajan K, Sali A. 4-Dimethylaminopyridine-catalyzed multi-component one-pot
411 reactions for the convenient synthesis of spiro[indoline-3,4'-pyrano[2,3-c]pyrazole]
412 derivatives. *Tetrahedron*. 2014;70(2):484-9.

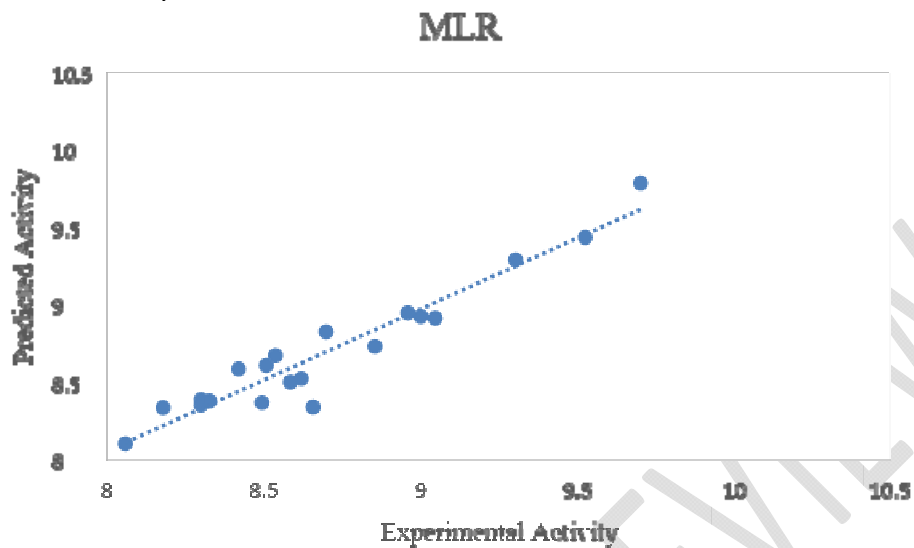
- 413 27. A. Fassihi, D. Abedi, L. Saghaie, R. Sabet, H. Fazeli, G. Bostaki, O. Deilami, H. Sadinpour,
414 Eur. J. Med. Chem. (2008), doi: 10.1016/j.ejmech.2008.10.022.
- 415 28. Sharaf MA, Illman DL, Kowalski BR. Chemometrics. New York: Wiley. 1986;332.
- 416 29. Olah, M.; Bologna, C.; Oprea, T.I. An Automated PLS Search for Biologically Relevant
417 QSAR Descriptors. J. Comput. Aided Mol. Des. 2004, 18, 437-449.
- 418 30. Mohajeri, A.; Hemmateenejad, B.; Mehdipour A.; Miri, R. Modeling Calcium Channel
419 Antagonistic Activity of Dihydropyridine Derivatives Using QTMS Indices Analyzed by
420 GA-PLS and PC-GA-PLS. J. Mol. Graph. Model. 2008, 26, 1057-1065
- 421 31. Brereton R. Chemometrics Data Analysis for the Laboratory and Chemical Plant. Wiley.
422 2004:47–54.
- 423 32. Liu. Y and Liu.T, Recent in non-peptidomimetic dipeptidyl peptidase 4 inhibitors:
424 Medicinal chemistry and preclinical aspects. Current Medicinal Chemistry, 2012, 19, 3982-
425 3999.

426
427
428
429 **Figure 1.** Principal component analysis diagram for detection of outlier data
430

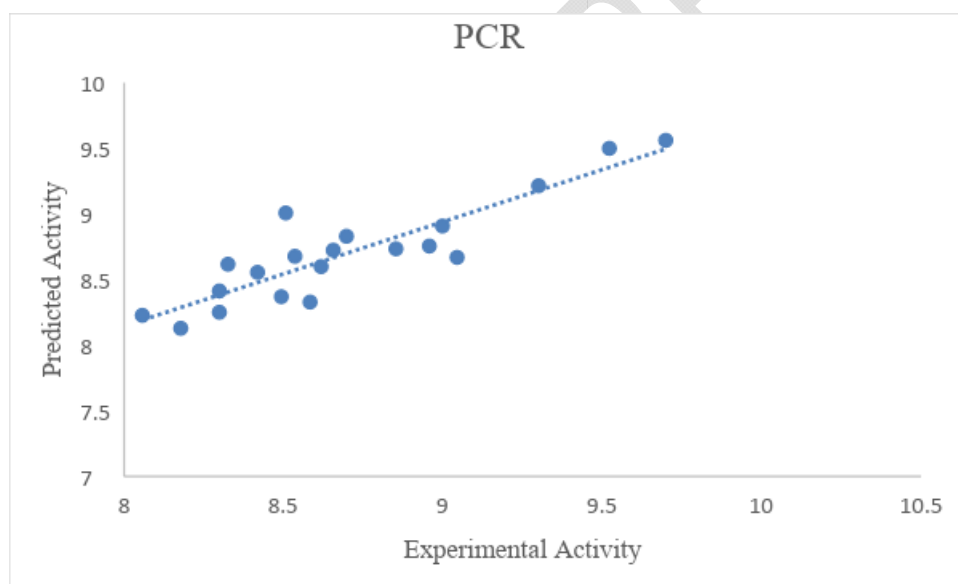


431
432

433 **Figure 2.** Plots of the cross-validated predicted activity against the experimental activity for the
434 QSAR models obtained by different methods.

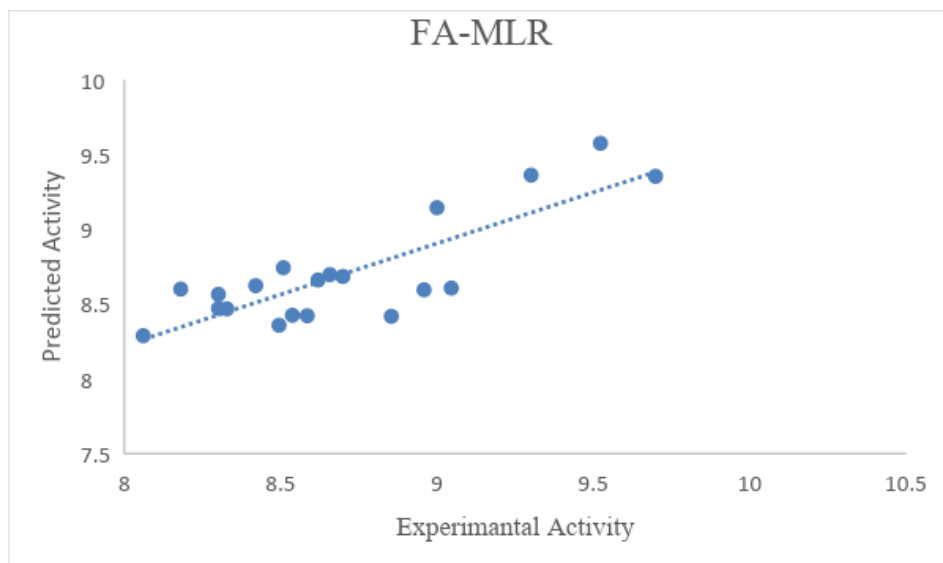


435



436

437



438

439

440

441

442

443

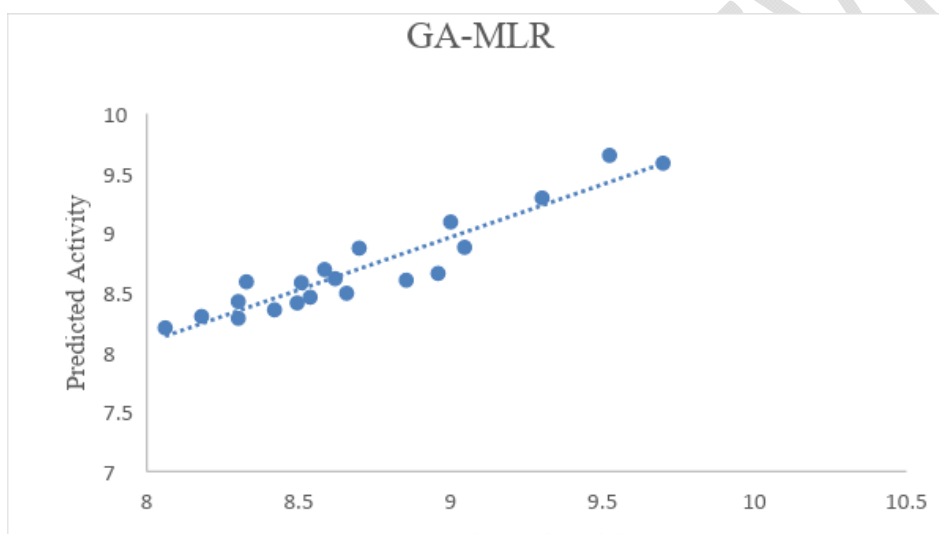
444

445

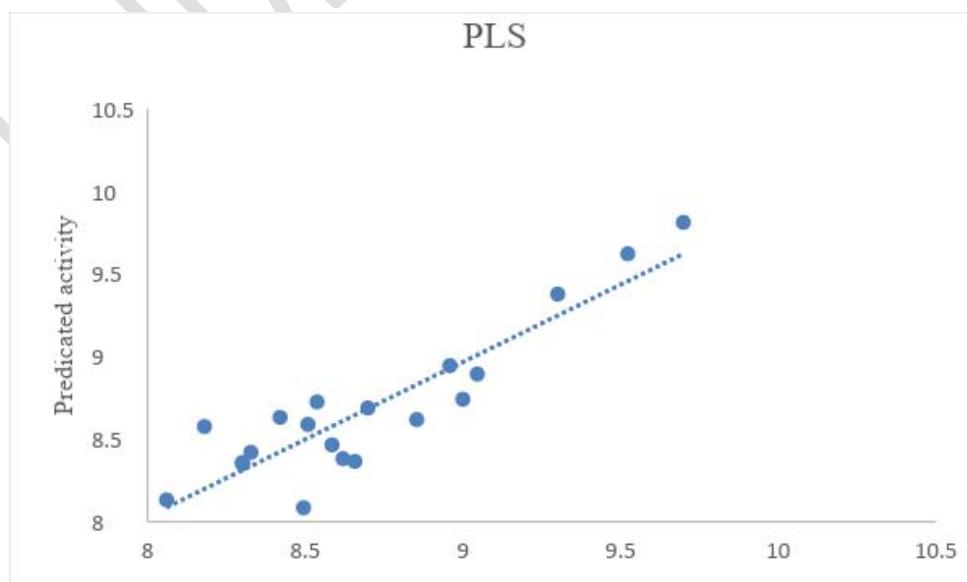
446

447

448



449



458

459

460

UNDER PEER REVIEW

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

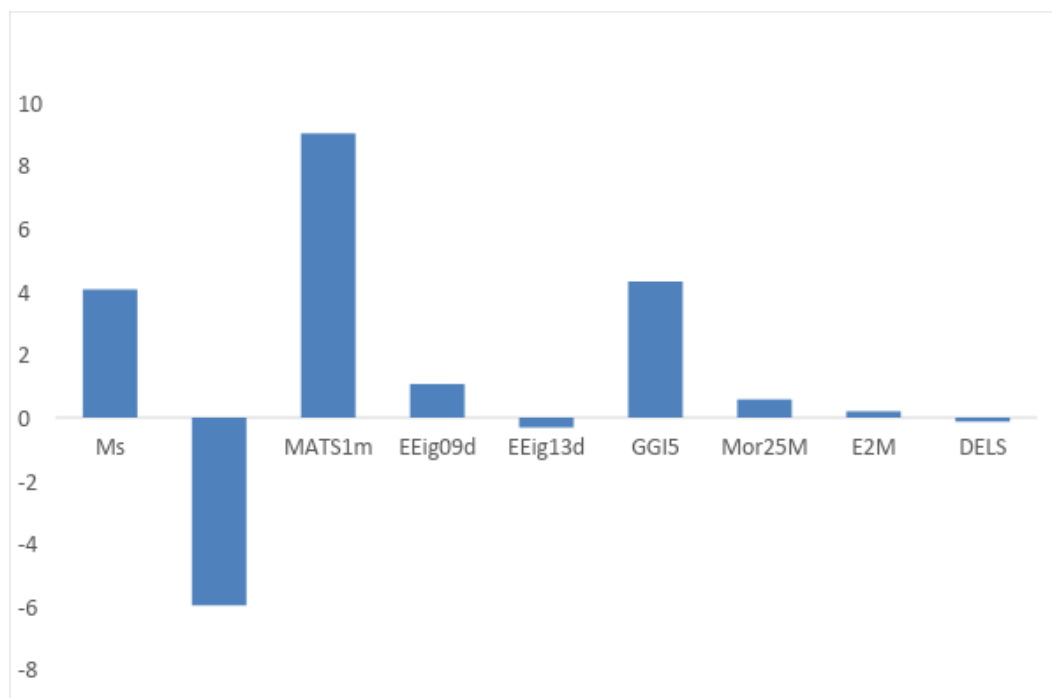
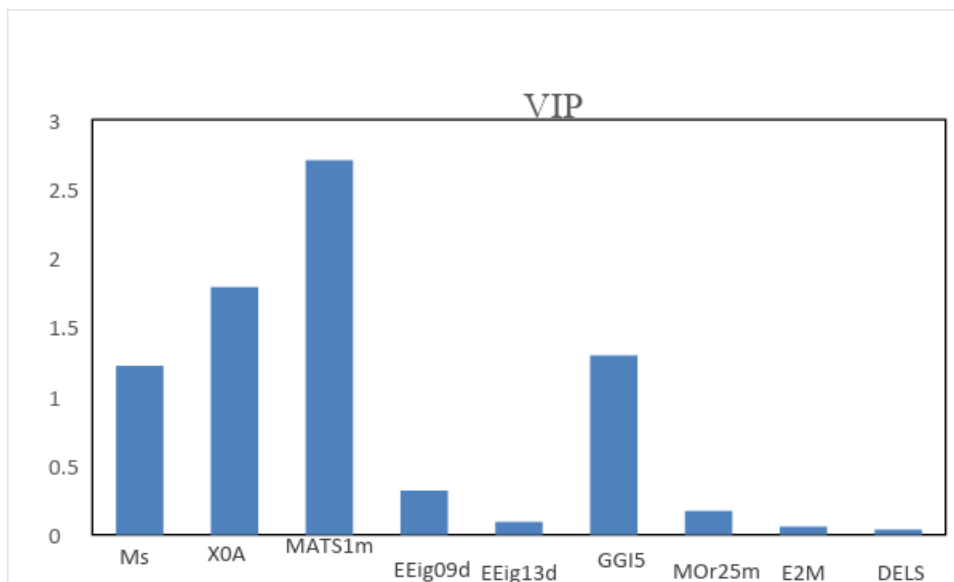
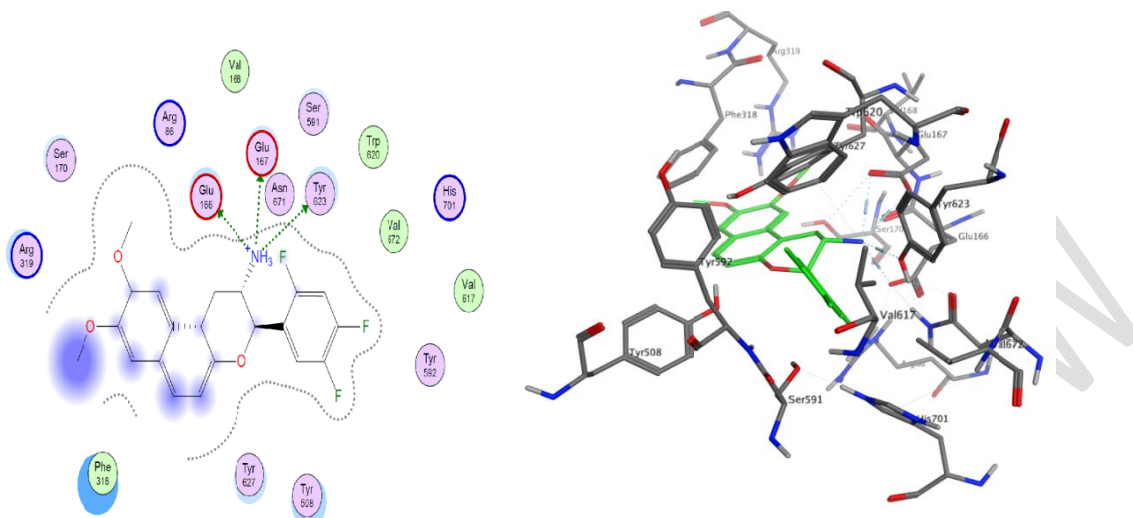


Figure 3. Plots of the cross-validated predicted activity against the experimental activity for the QSAR models obtained by GA-PLS methods.

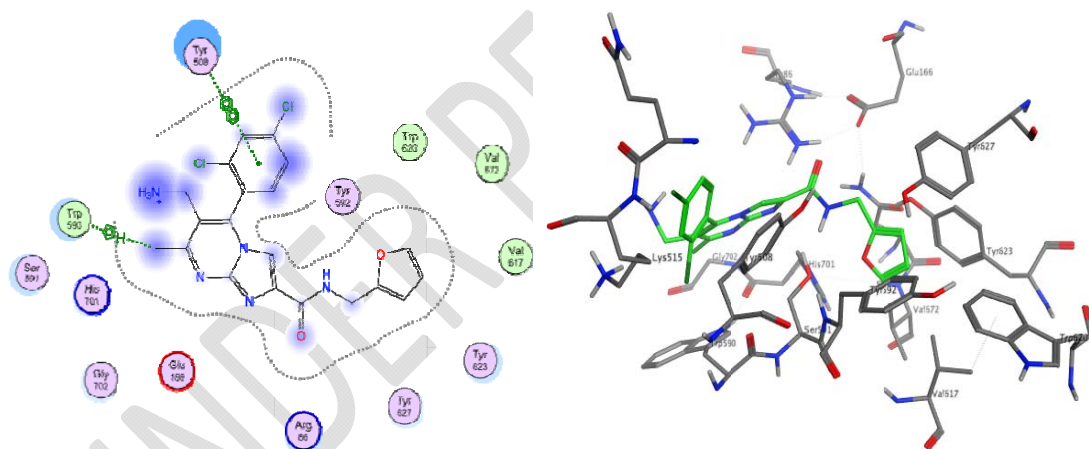


481
482 **Figure 4.** Plot of variables important in projection (VIP) for the descriptors used in GAPLS
483 model.
484

485 **Figure 5.** Interactions of A) HL1 and B) compound 25 with the residues in the binding site of DPP4 (5j3j)
486 receptor.
487



488
489
490
491
492
493



494
495
496
497
498
499

500 **Table1.** Chemical structure of imidazolopyrimidine amides analogues used and their
501 experimental and cross validated-predicted activity by (GA-PLS) for DPP4 inhibitory and their
502 docking bonding energies.

503

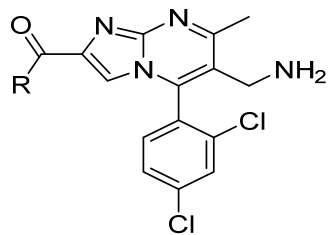
504

505

506

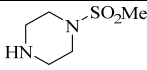
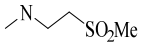
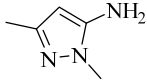
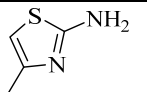
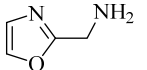
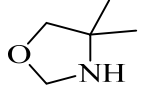
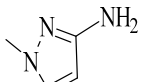
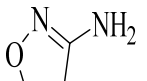
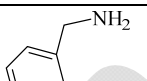
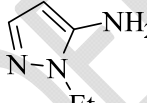
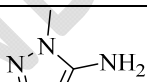
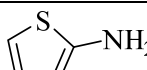
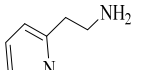
507

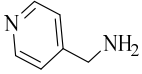
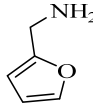
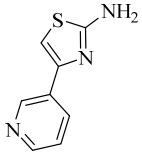
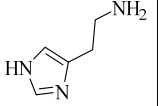
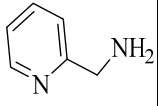
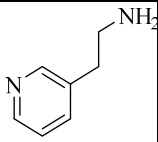
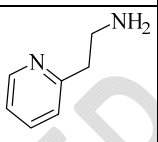
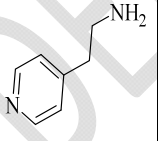
508



1-31

NO	R	Exp.pIC ₅₀	Pred. pIC ₅₀ by GA-PLS	Binding Energy (kcal/mol)
1*	OEt	9.39	-----	-----
2		8.6	8.38	-8.1
3		8.5	8.7	-8.7
4		8.6	8.3	-7.9
5**		8.5	8.65	-8.6
6		8.3	8.3	-8.1
7		8.06	8.1	-8.2
8		9	8.7	-8.1
9		8.5	8.5	-8
10		9.69	9.8	-8.2

11		9.3	9.3	-8.2
12		9.5	9.6	-7.9
13		9.04	8.9	-8.9
14		8.18	8.5	-8.4
15		8.69	8.7	-9
16*		-----	-----	-----
17**		8.7	8.37	-8.5
18		8.40	8.60	-9.3
19**		8.58	8.3	-8.9
20		8.95	8.9	-8.6
21**		8.69	8.68	-8.6
22		8.49	8.48	-8.7
23**		9.15	9.1	-8.9

24		8.58	8.4	-8.4
25**		8.39	8.4	-9.4
26		8.3	8.35	-9.2
27		8.32	8.4	-8.3
28**		8.26	8.28	-9.4
29**		8.8	8.72	-8.8
30		8.8	8.6	-8.9
31**		8.8	8.73	-8.3

509 *: outlier data

510 **: molecules as test set

511

512

513

514 **Table 2.** Brief name of molecular descriptors was used in the models.

Descriptor type	descriptors	Brief description
Constitutional	Ms.	Mean electropological state
Topological	Jhetm	Balaban-type index from mass weighted distance matrix
	DELS	Molecular electropological variation
Connectivity indices	X0A	Average connectivity index chi-0
2D-autocorrelation	MATS1m	Moran autocorrelation – lag1/weighted by atomic Masses
Edge adjacency indices	EEig09d	Eigen values 09 from edge adj. matrix weighted by dipole moment
	EEig13d	Eigen values 13 from edge adj. matrix weighted by dipole moment
Burden Eigenvalues	BELm6	Lowest eigenvalue n.6 of burden matrix/weighted by atomic masses
Topological charge indices	GGI5	Topological charge index of order 5
	GGI7	Topological charge index of order 7
3-D Morse Descriptors	Mor27u	3D-MoRSE-signal 27/unweighted
	Mor25m	3D-MoRSE-signal 25/weighted by atomic masses
WHIM Descriptors	E2m	2 nd component accessibility directional WHIM index/weighted by atomic masses

515

516

517

518

519

520 **Table 3.** The results of different QSAR model analysis

Models	Equation	N	R ²	Q ²	F	SE	R ² _p
MLR	PIC ₅₀ =9.508 MATS1m (±2.252) +4.286 Ms (±0.78) +4.319 GGI5 (±0.816)-0.105 DELS (±0.028) +0.538 MOR25m (±0.182)-3.903(±1.868)	29	0.91	0.84	31.7	0.14	0.92
PCR	PIC ₅₀ = 0.245 PC1(±0.243) +0.13 PC3 (±0.043) +0.129 PC2(±0.043)-0.121 PC7(±0.043) + 8.695(±0.042)	29	0.77	0.75	14.6	0.22	0.83
FA-MLR	PIC ₅₀ =11.953 MATS1m(±2.8) +2.65 Ms (±0.83) +2.61(±1.96)	29	0.67	0.53	13.2	0.12	0.63
GA-MLR	PIC ₅₀ =2.072GGI7((±0.676)+4.427Ms((±0.678)+ 8.047BELm6(±1.753)-0.453Mor27u(±0.187)-14.411(±3.275)	29	0.94	0.88	28.5	0.16	0.91
GA-PLS	-----	29	0.94	0.80	-----	0.49	0.95

521

522

523

524 **Table 4.** Correlation coefficient (R²) matrix for descriptors represented in multiple linear
525 regression

526

	MATS1m	Ms	GGI5	DELS	Mor25m
MATS1M	1	0.413	0.649	0.712	0.077
MS		1	0.345	0.668	0.250
GGI5			1	0.859	-0.125
DELS				1	0.079
Mor25M					1

527

528

529

UNDER PEER REVIEW

530 **Table 5.** Numerical values of factor loading numbers 1–7 for some descriptors after VARIMAX
 531 rotation (against DPP4 inhibitory activity).

	Component						
	1	2	3	4	5	6	7
volume	.524	.262	.073	.125	.214	-.110	-.325
Ms	.252	.792	.340	.019	-.208	.198	-.119
nH	.472	-.700	.088	.075	.369	-.295	.113
STN	-.163	.087	-.754	-.030	.464	.225	.058
DELS	.753	.448	.131	.132	.283	.005	.031
X0A	.335	.038	.881	.115	.084	-.127	.040
IVDE	.122	.190	.905	.185	.145	-.012	.150
IC0	.202	.925	.173	.205	.007	.043	.079
MATS1m	.849	.054	.267	-.049	.046	.041	-.155
GATS6m	-.179	-.646	-.430	-.251	-.120	.299	-.231
EEig09d	.670	.296	.117	.081	.372	-.295	.163
EEig13d	.252	-.104	-.570	-.063	.548	-.358	.176
GGI5	.598	.295	.270	.278	.508	-.093	.119
GGI4	-.054	.093	.319	.058	.785	.052	.239
JGT	.008	.319	.842	.136	.283	.035	.135
RDF015m	.874	-.157	-.142	-.121	.021	.311	.123
Mor04m	-.029	-.482	-.330	.096	-.442	.190	.180
Mor25m	.053	-.006	-.130	-.055	-.079	-.003	-.885
Mor19m	.243	.052	-.009	-.043	-.176	.905	.179
Mor18p	.383	-.212	.282	.034	-.226	-.699	.321
E2m	.822	.165	.140	.097	-.144	-.040	-.068
HATS3e	-.304	.081	.133	.023	-.752	.080	.171
R4e	.559	-.465	.027	.139	.355	-.399	.210
ALOGP2	-.075	-.110	-.117	-.975	.006	.043	-.062
TE1	.321	.162	.195	.900	.056	-.012	-.012
TPSA(Tot)	.206	.817	-.099	.148	.346	.065	-.014
F03[C-O]	.175	-.031	-.064	-.976	-.001	.027	-.011

532

533

534

535 **Table 6.** Leverage (h) of the external test set molecules for different models. The last row (h^*) is
536 the warning leverage.

Molecular no	MLR	PCR	FA-MLR	GA-MLR	GA-PLS
5	0.13	0.07	0.049	0.37	1.09
17	0.26	0.10	0.052	0.26	0.3
19	0.13	0.20	0.048	0.20	0.65
21	0.23	0.15	0.052	0.19	0.30
23	0.089	0.25	0.047	0.32	1.0
25	0.076	0.065	0.052	0.137	0.37
28	0.129	0.23	0.047	0.114	0.52
29	0.24	0.095	0.048	0.100	0.34
31	0.3	0.105	0.048	0.10	0.48
h^*	0.71	0.6	0.3	0.6	1.35

537

538

539