

Decoding the inhibition selectivity and physicochemical multiverse of beta-carbonic anhydrase inhibitors through random forest and neural network-assisted QSAR/QSAAR modeling: A novel approach for the rational design of new anti-tuberculosis drugs

- **Principle Investigator:** Dr. Ashok Aspatwar
- **ML Model Generation and Optimization/Software Development/Chemical Data Curation/ Molecular Modeling Studies:** Ratul Bhowmik
- **Co-authors:** Ajay Manaithiya, Rajarshi Ray
- **Affiliation:** Tampere University, Finland

-5.000

5.000

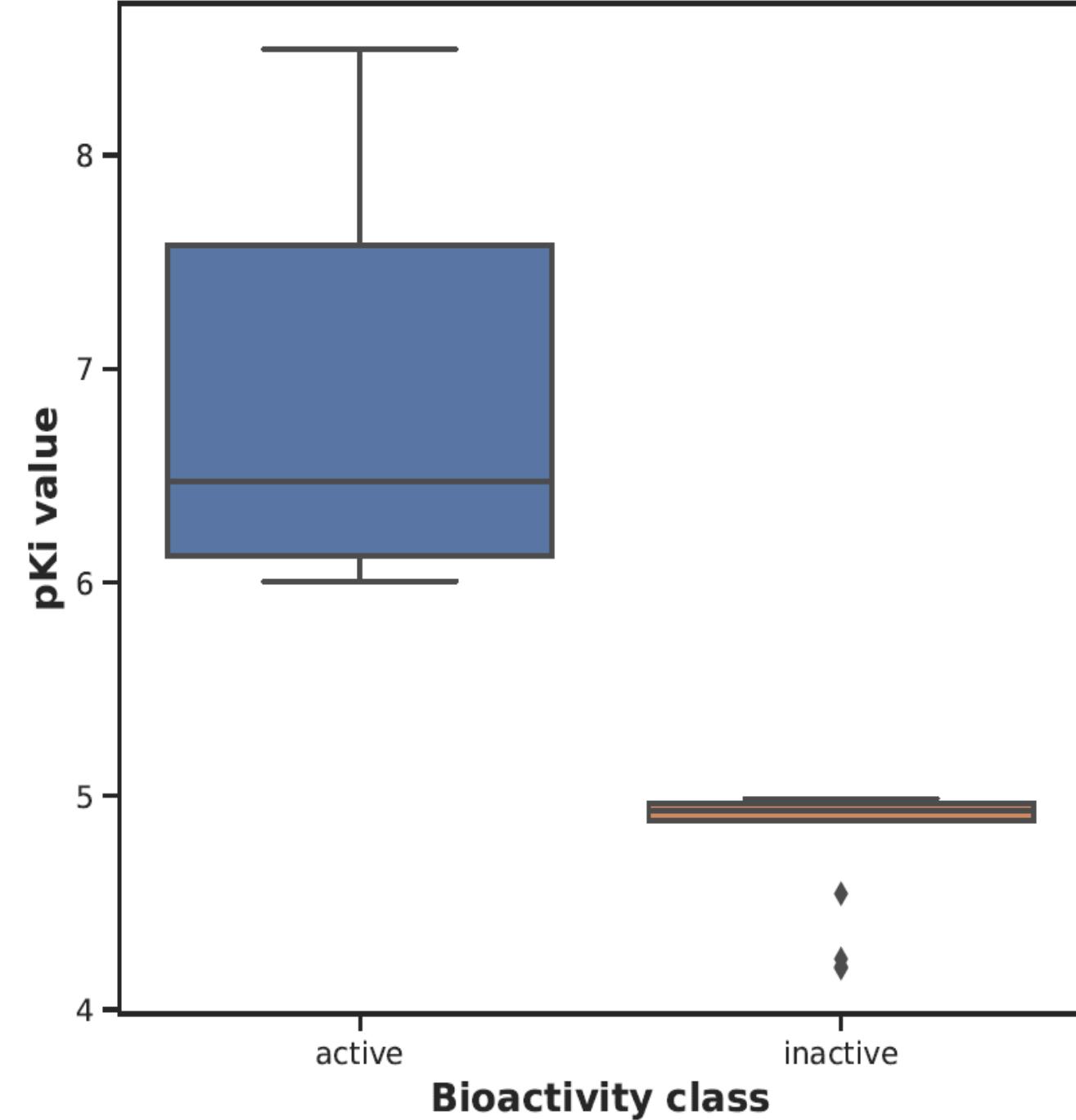
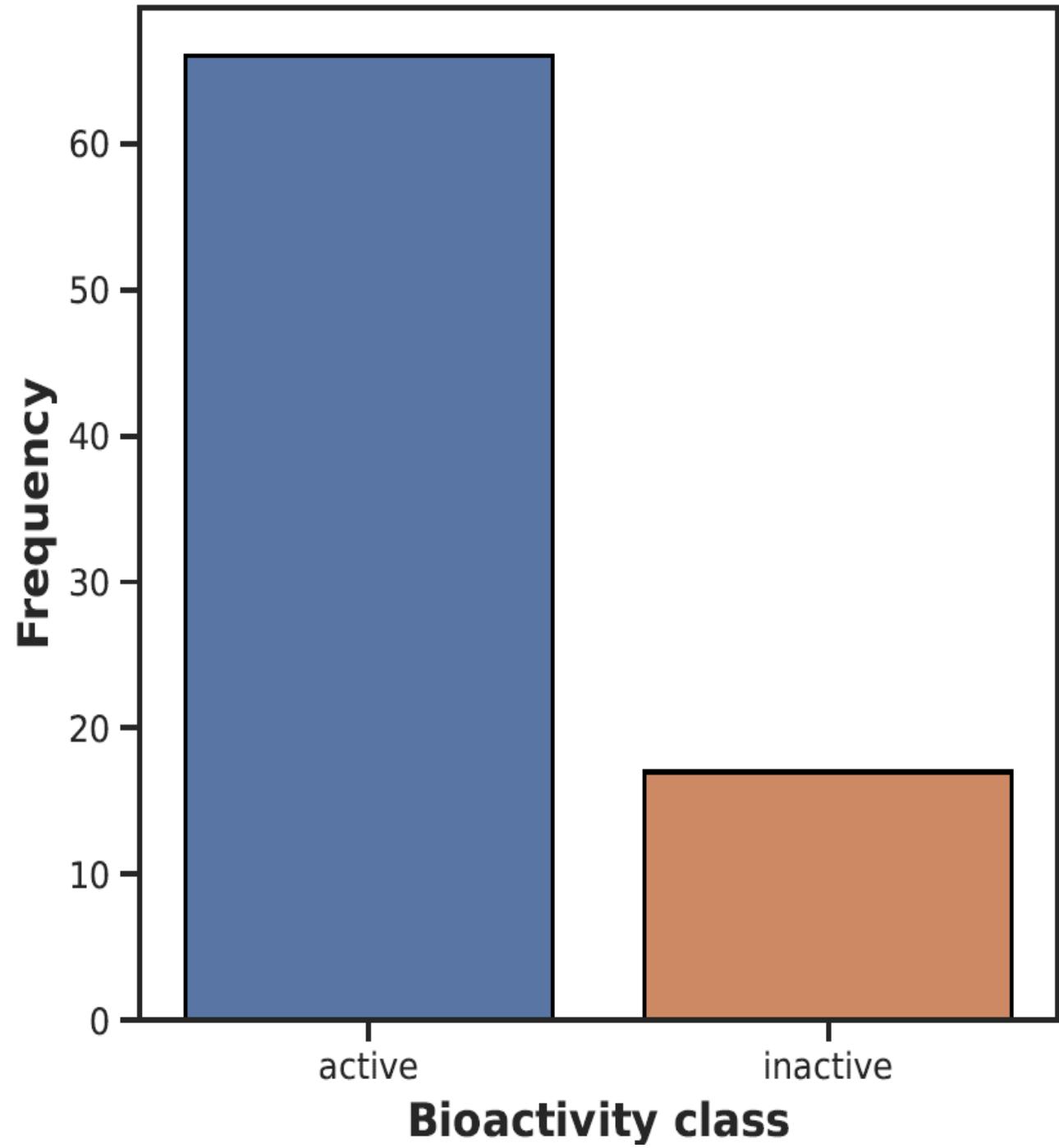
PRIMARY OBJECTIVES

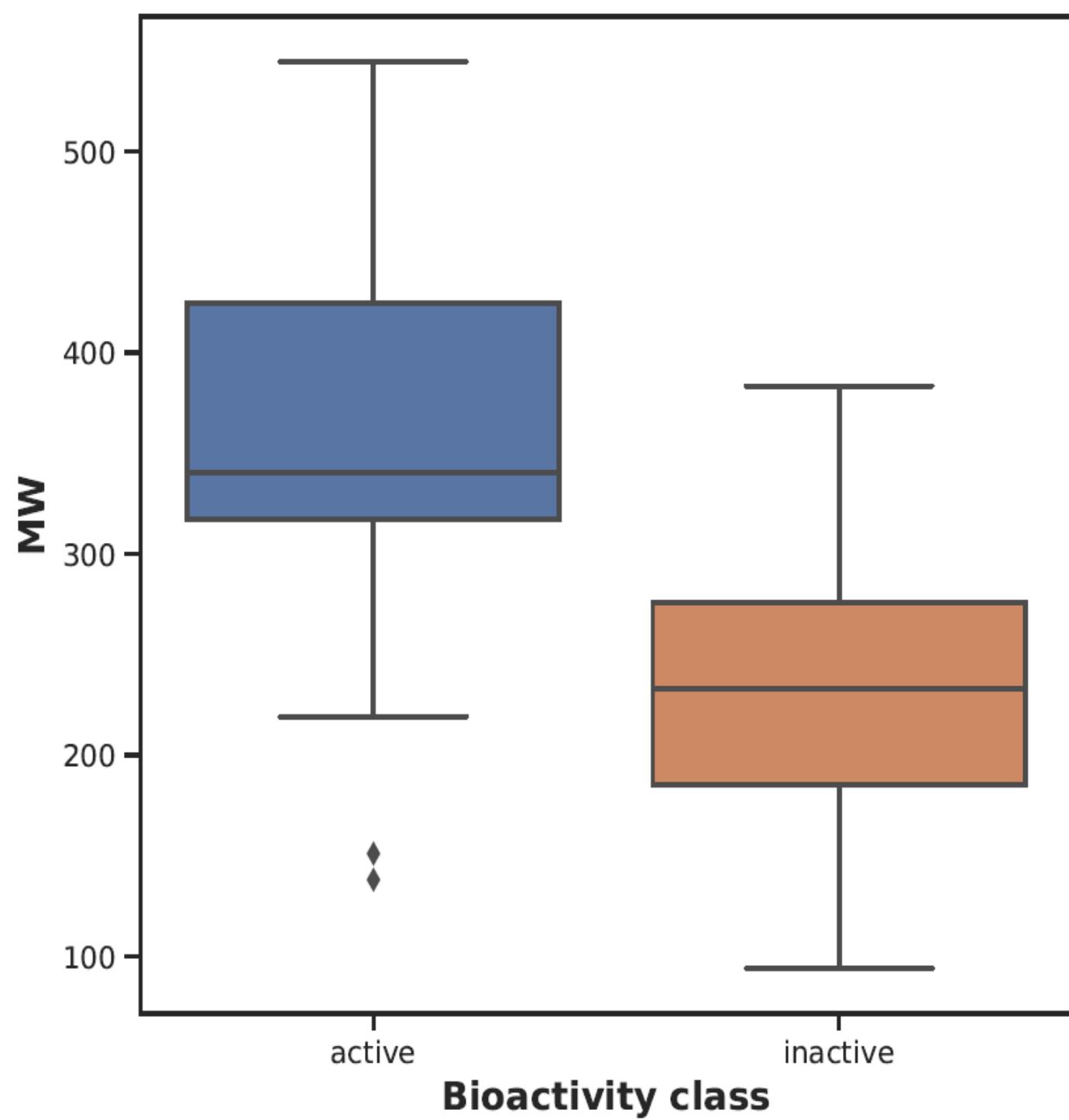
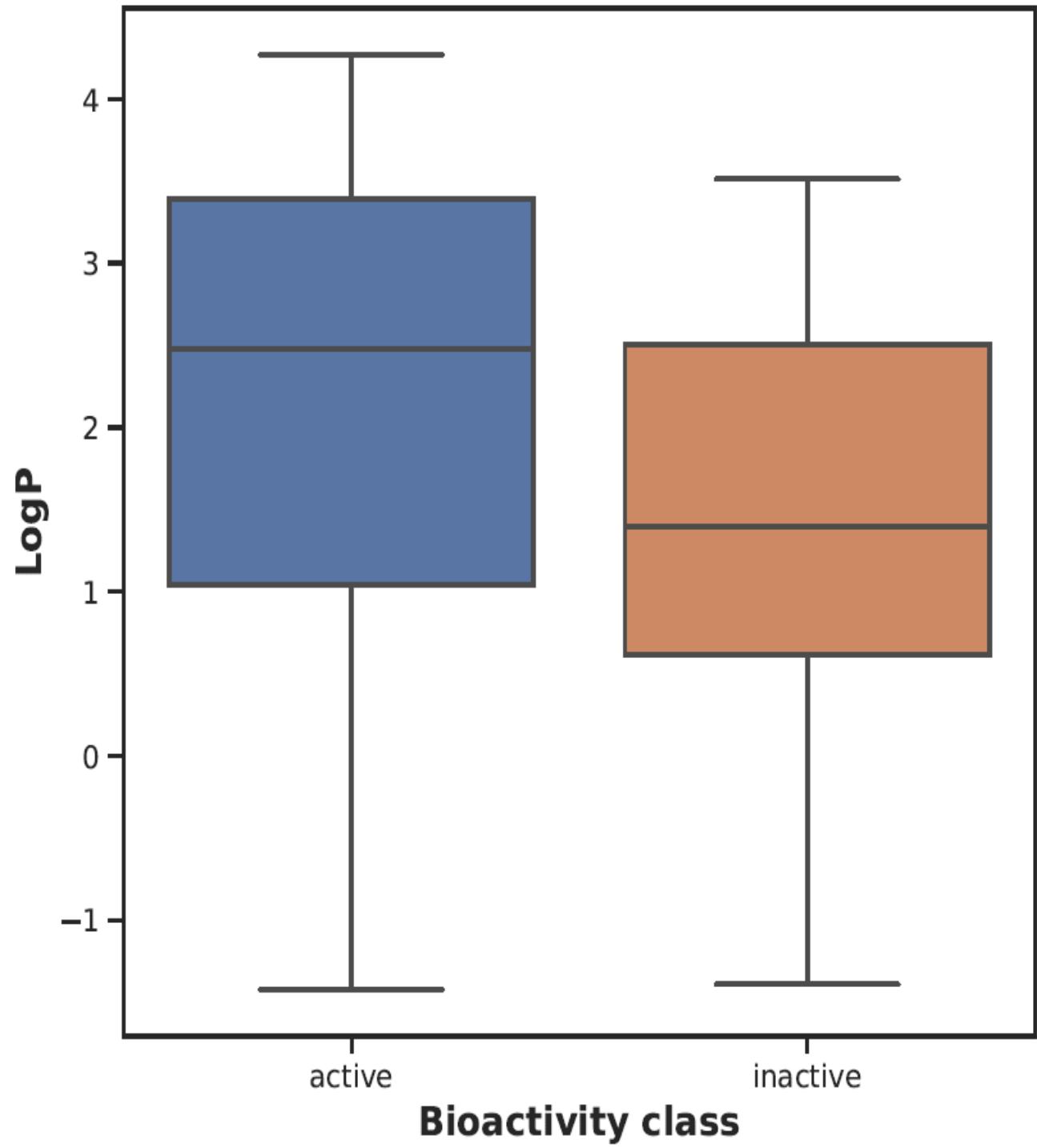
- Generation of selective robust multi-molecular signature implemented machine learning-assisted quantitative structural activity relationship (**ML-QSAR**) models for predicting bioactivity (**K_i**) against **MtbCA1** and **MtbCA2**
- Structural interpretation of active and inactive molecules concerning **MtbCA1** and **MtbCA2** inhibition through variance importance plot, correlation matrix, and feature plot analysis
- Development of a web application using the generated ML-QSAR models
- Generation of neural network-assisted quantitative structural activity-activity relationship (**ML-QSAAR**) models to infer the structural-functional correlation between **MtbCA1** and **MtbCA2** inhibitor's binding mechanism
- Validation of key molecular signatures extracted from **ML-QSAR** and **ML-QSAAR** analysis through molecular docking approach

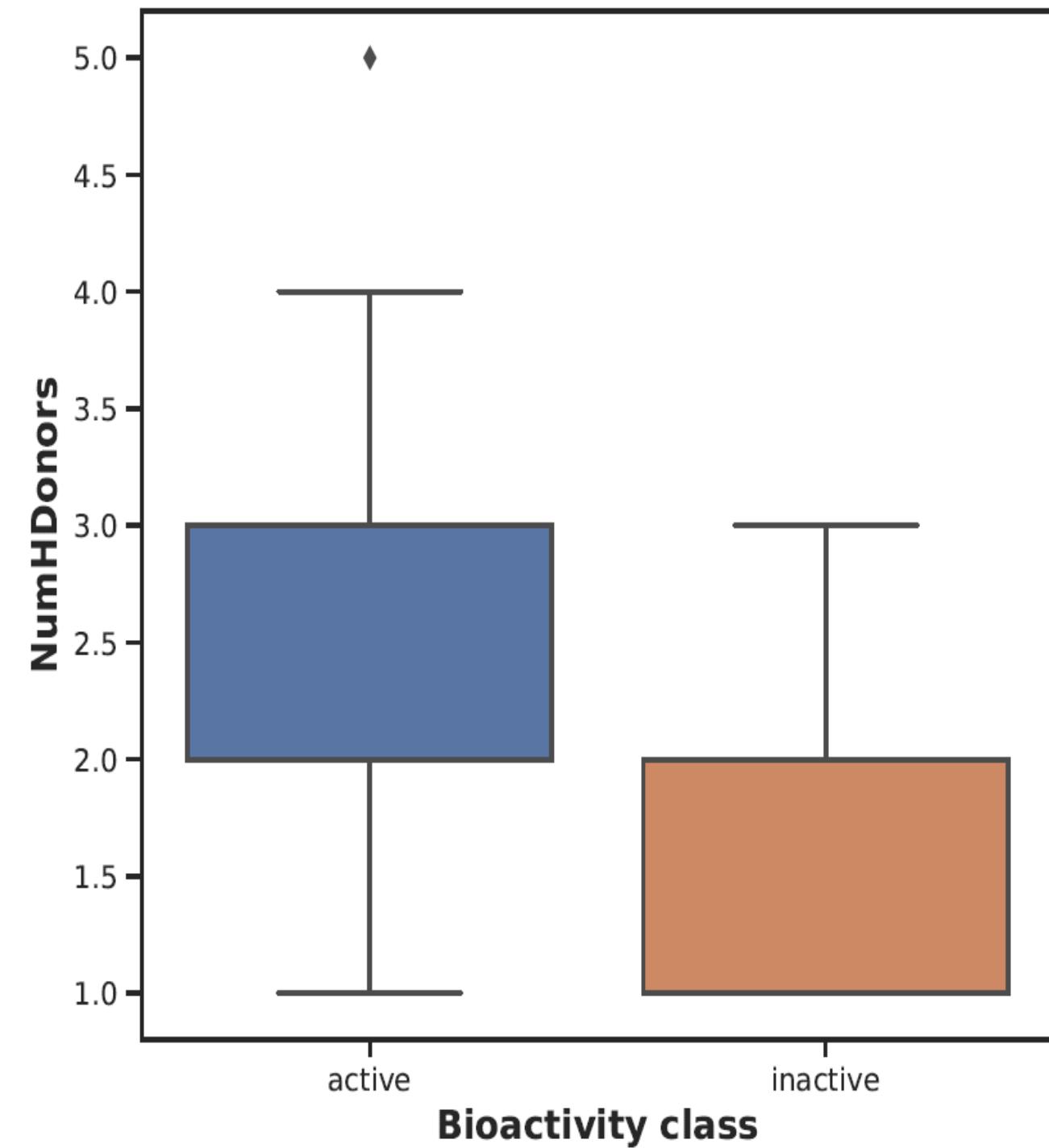
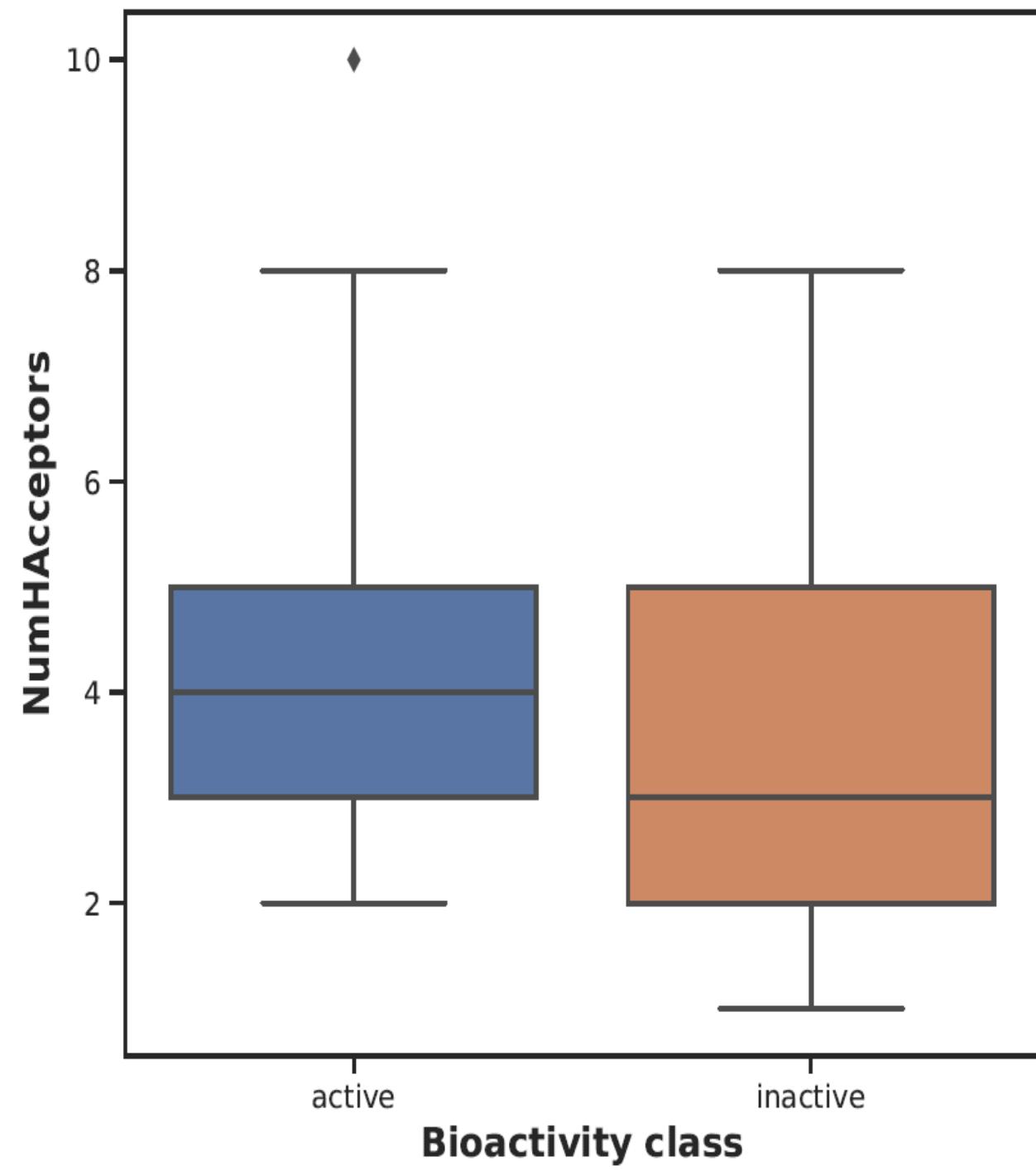
MtbCA1 Exploratory Data Analysis

-5.000

5.000



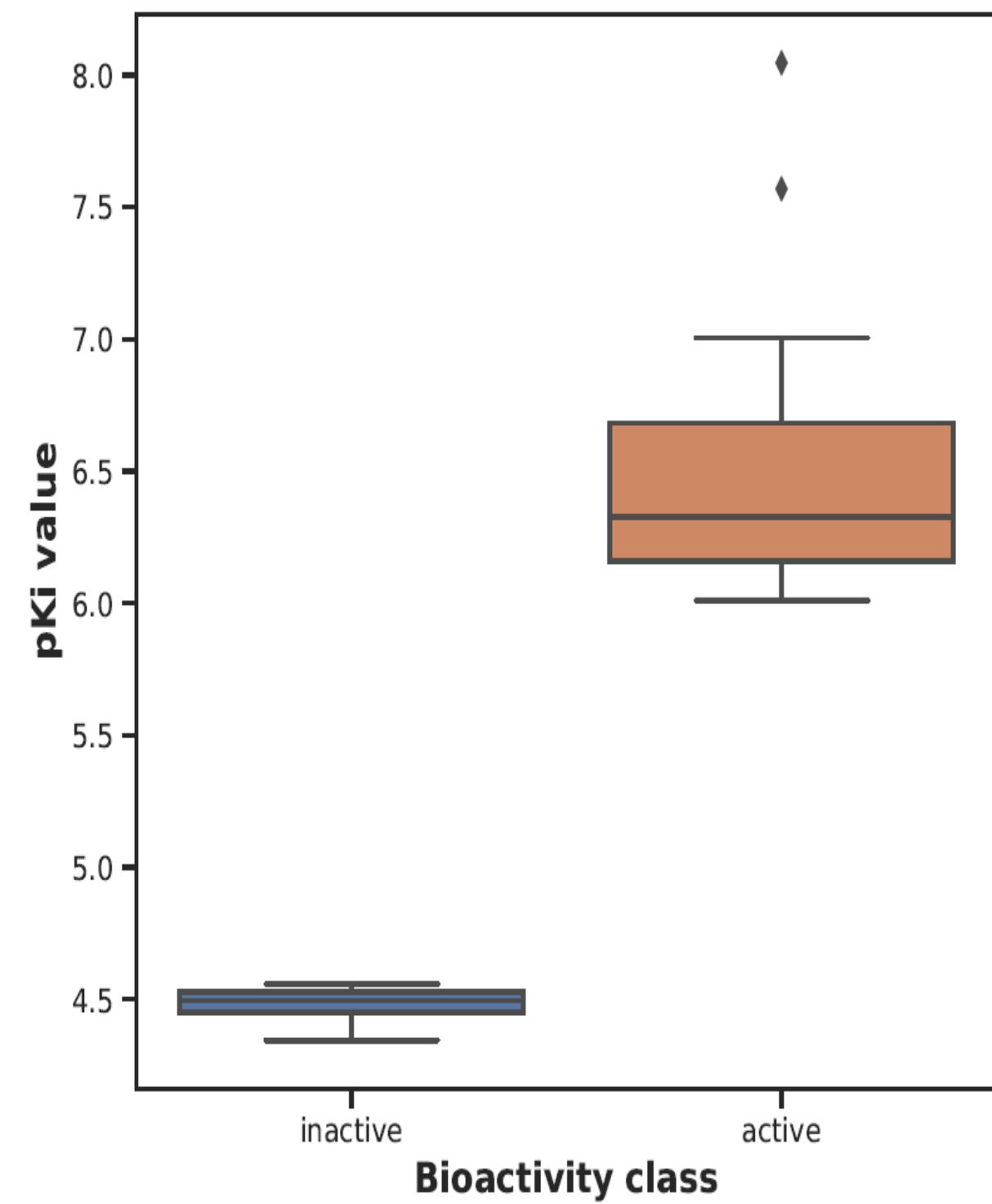
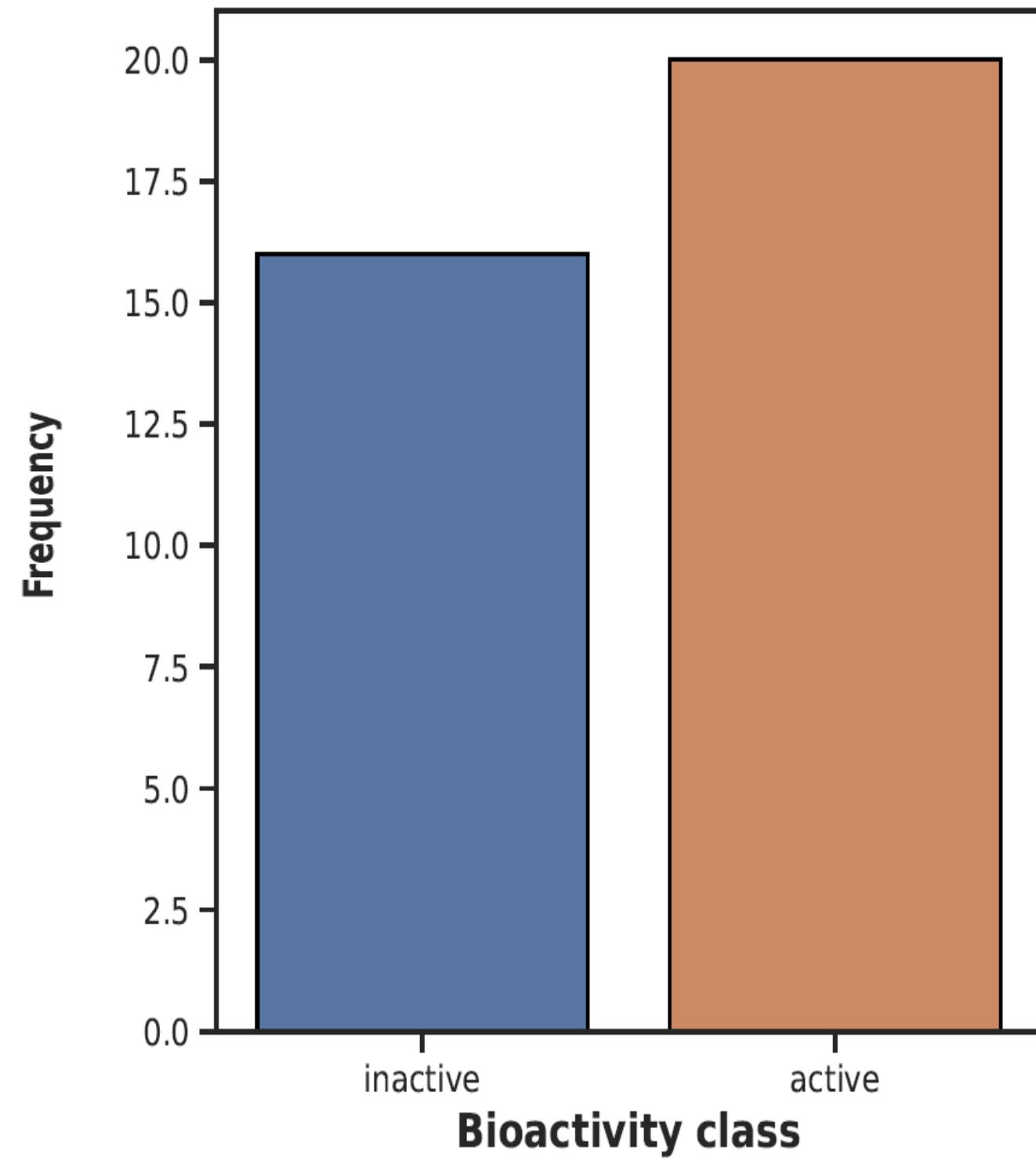


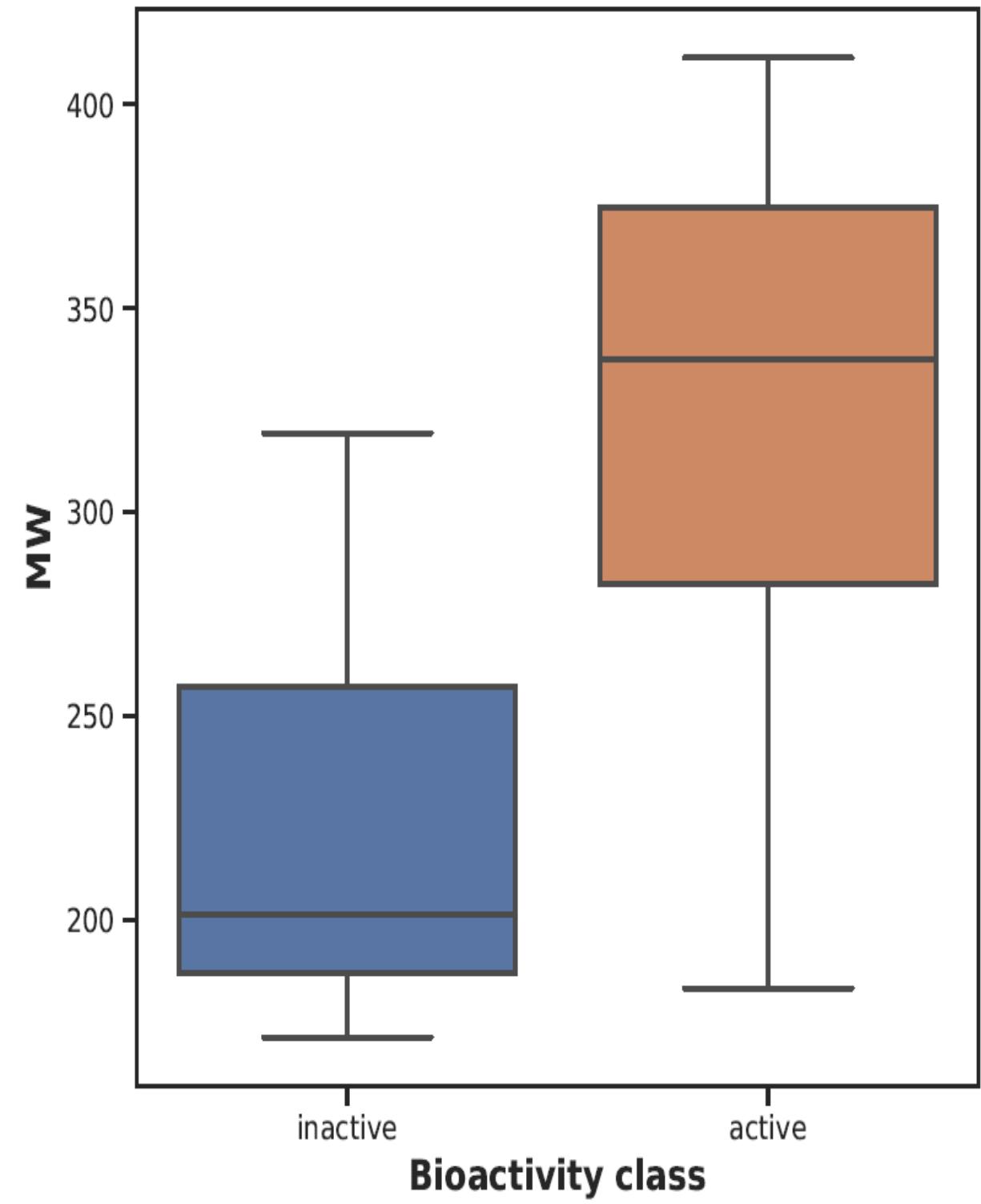
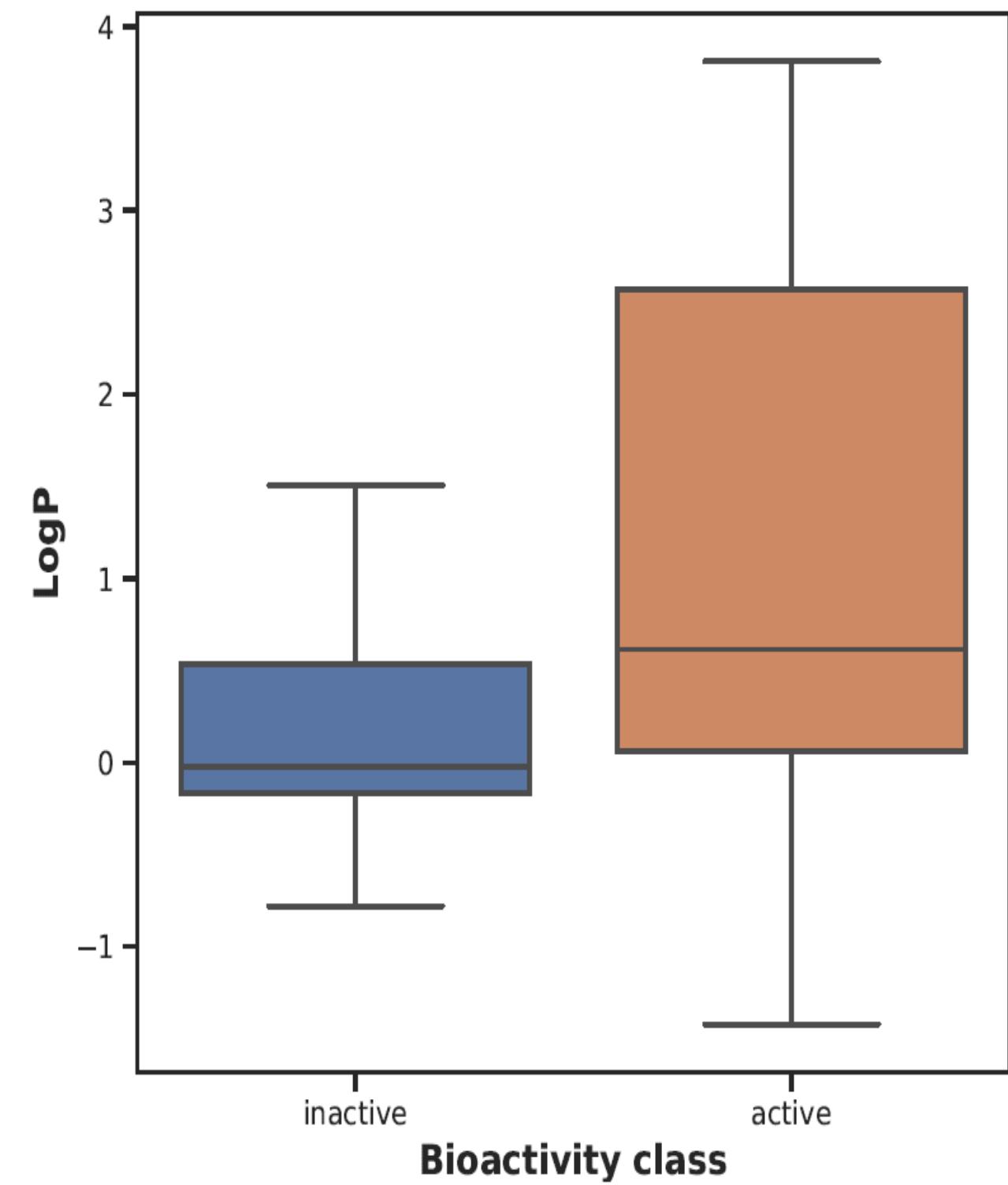


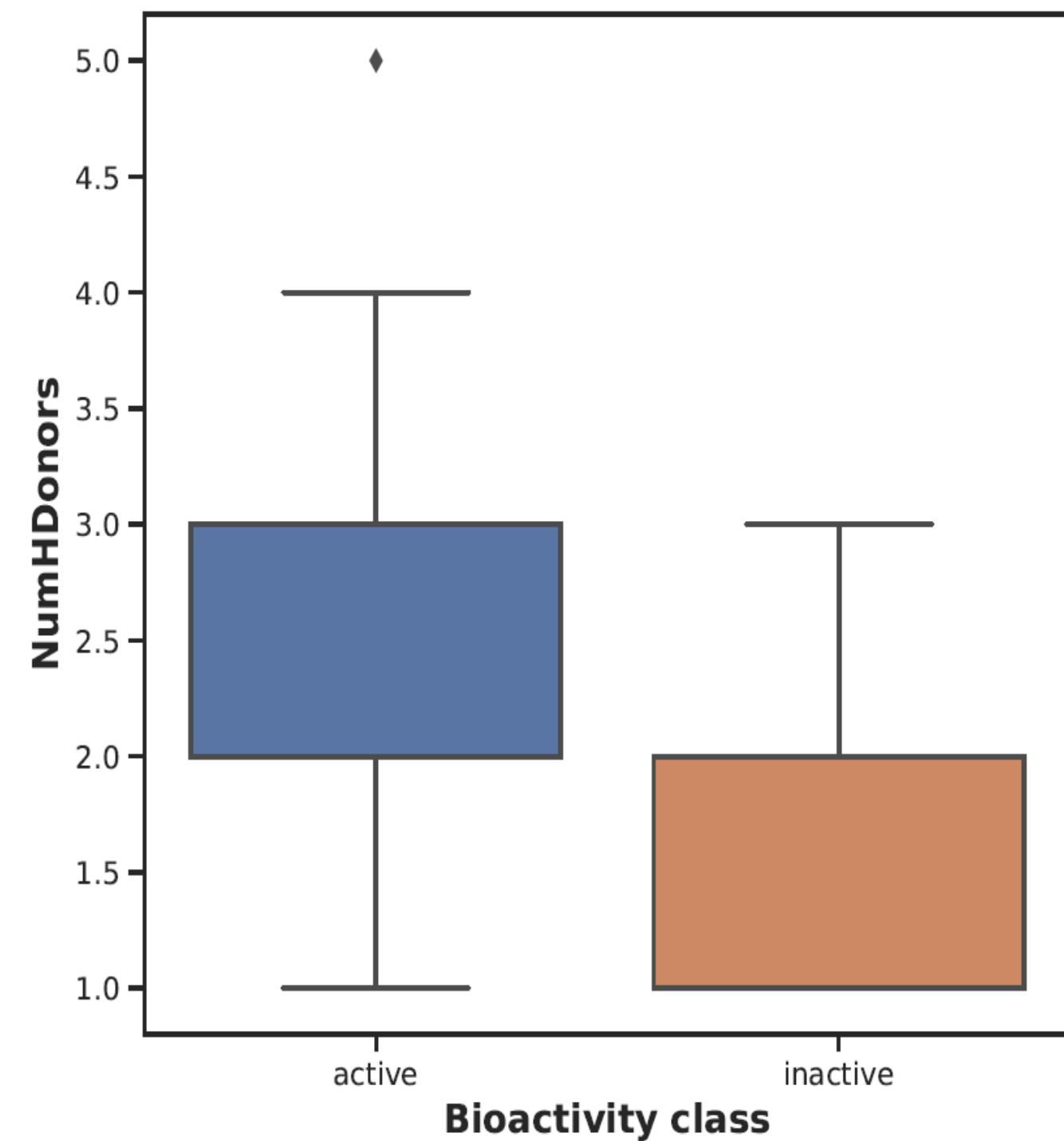
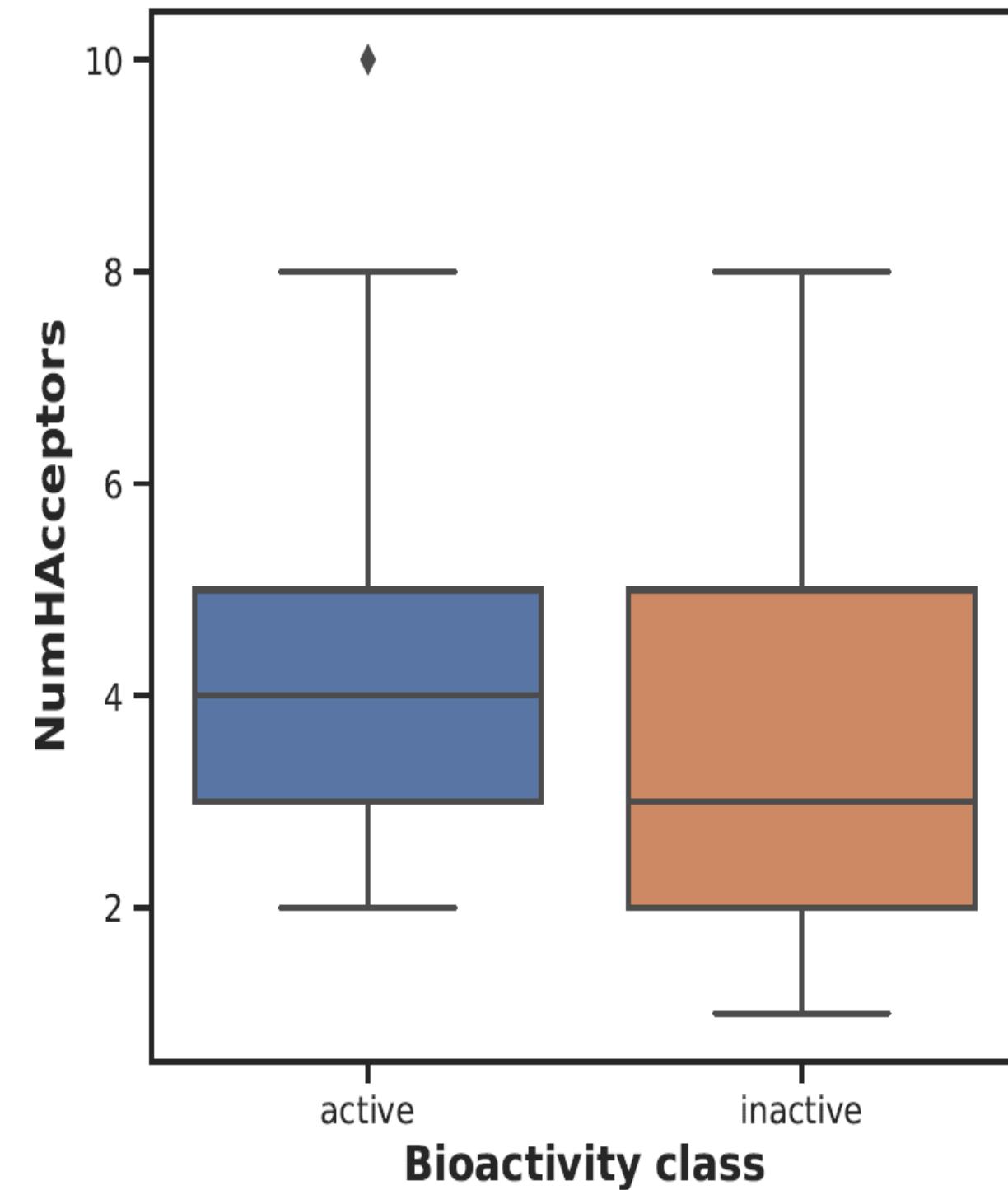
MtbCA2 Exploratory Data Analysis

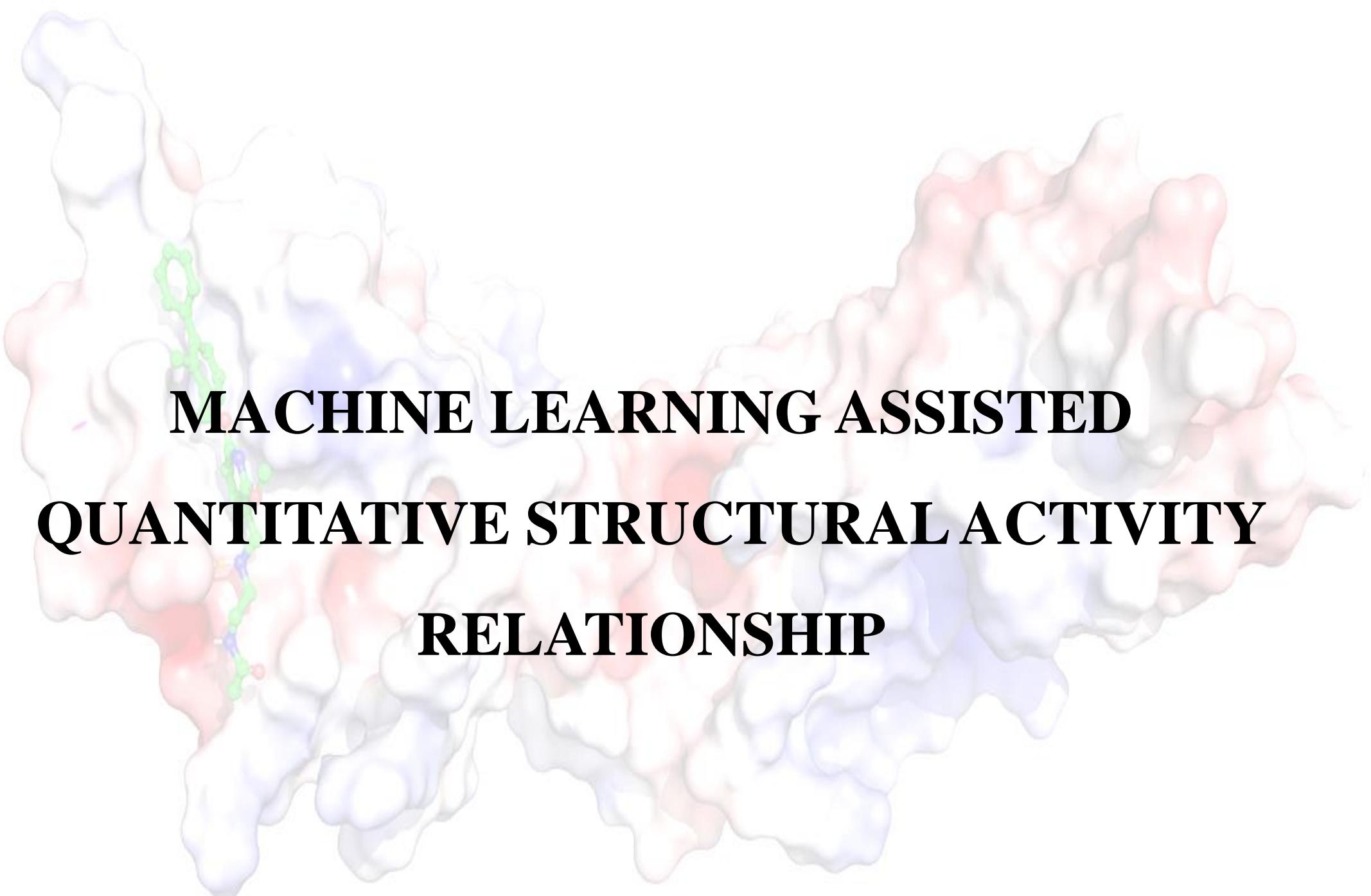
-5.000

5.000









MACHINE LEARNING ASSISTED QUANTITATIVE STRUCTURAL ACTIVITY RELATIONSHIP

-5.000

5.000

MtbCA1 Prediction Model

-5.000

5.000

ML-QSAR model for bioactivity prediction of *MtbCA1* inhibitors using substructure fingerprints



Curation of *MtbCA1* inhibitors with Ki value

Total inhibitors: 124

Substructure fingerprints: 307

Prediction Parameter: Bioactivity Class ('pKi')

Dataset Splitting into training and test set

Elimination of highly correlated and constant fingerprints

Final number of fingerprints(18), number of molecules in training(99) and test(25)

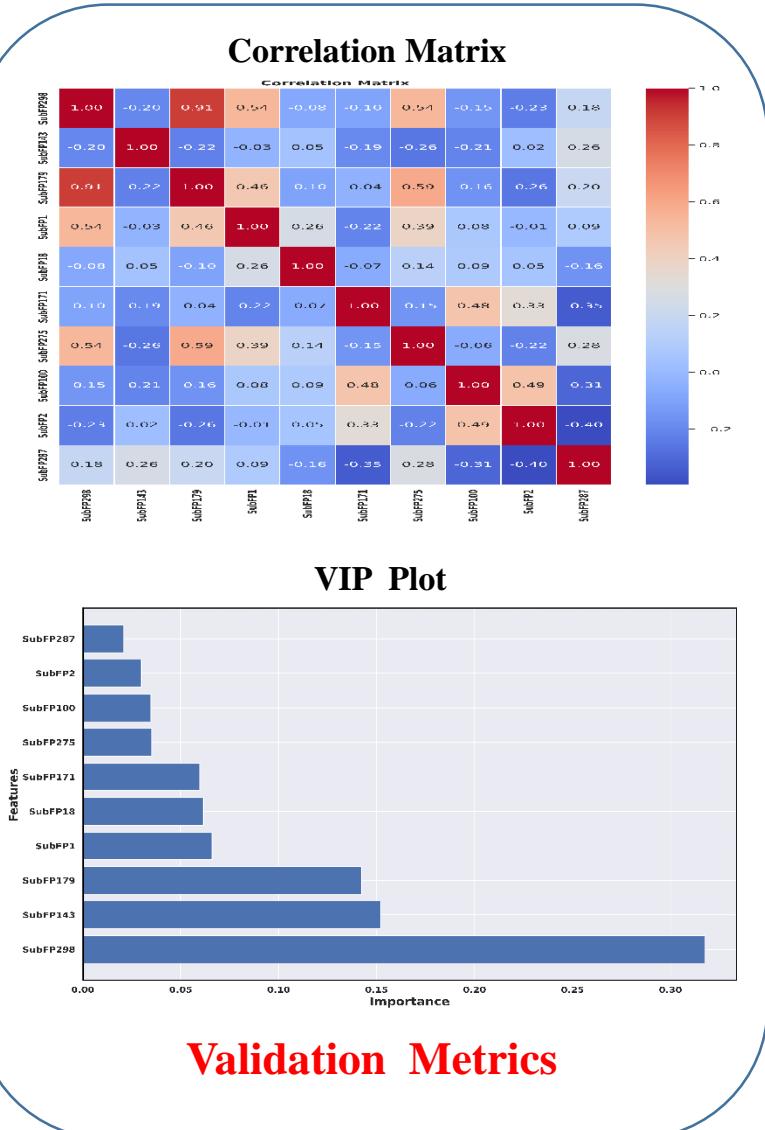
Random Forest Regressor

Initial Model Performance

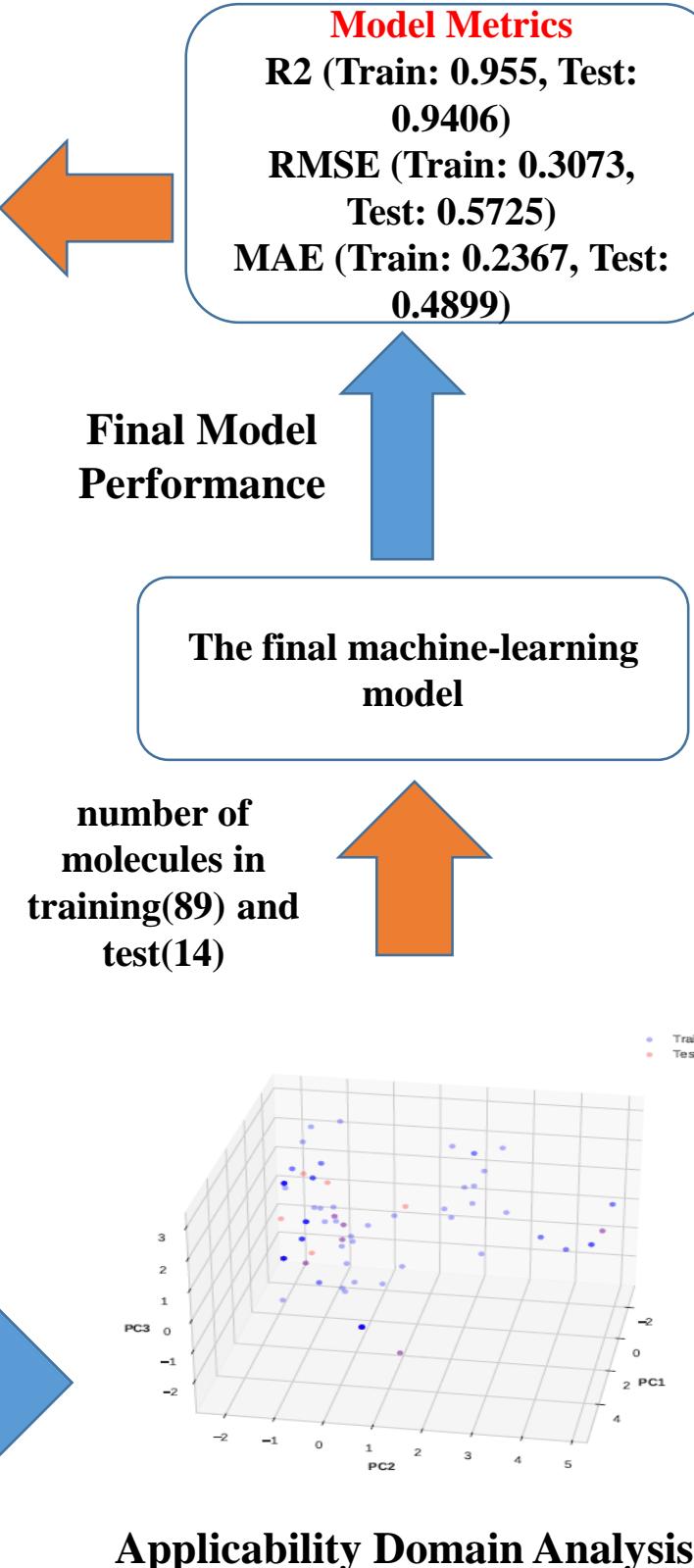
Model Metrics
 R2 (Train: 0.8853, Test: 0.7371)
 RMSE (Train: 0.5084, Test: 0.8106)
 MAE (Train: 0.3549, Test: 0.7171)

FP and FN Removal

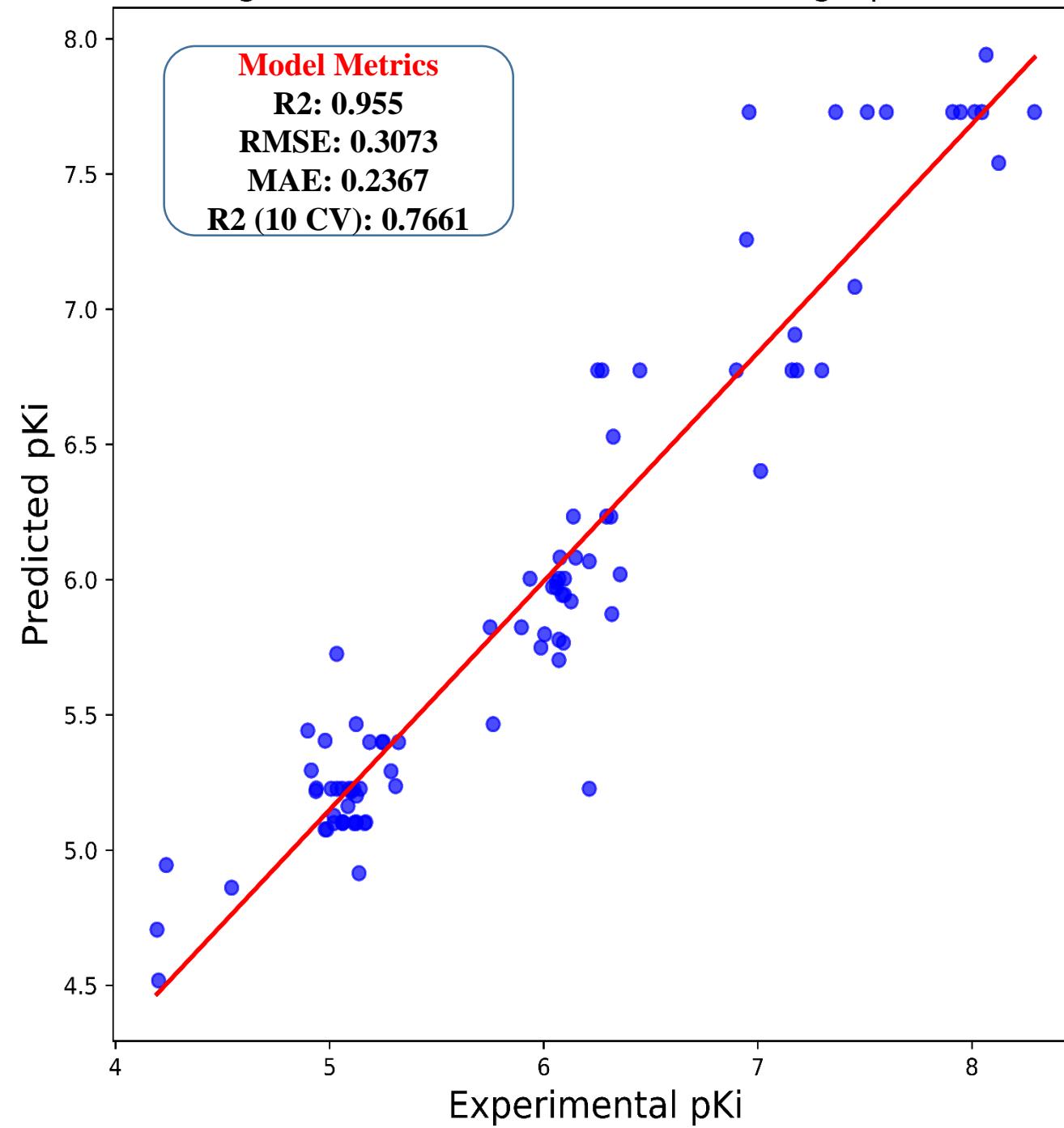
Mention molecules removed (21)



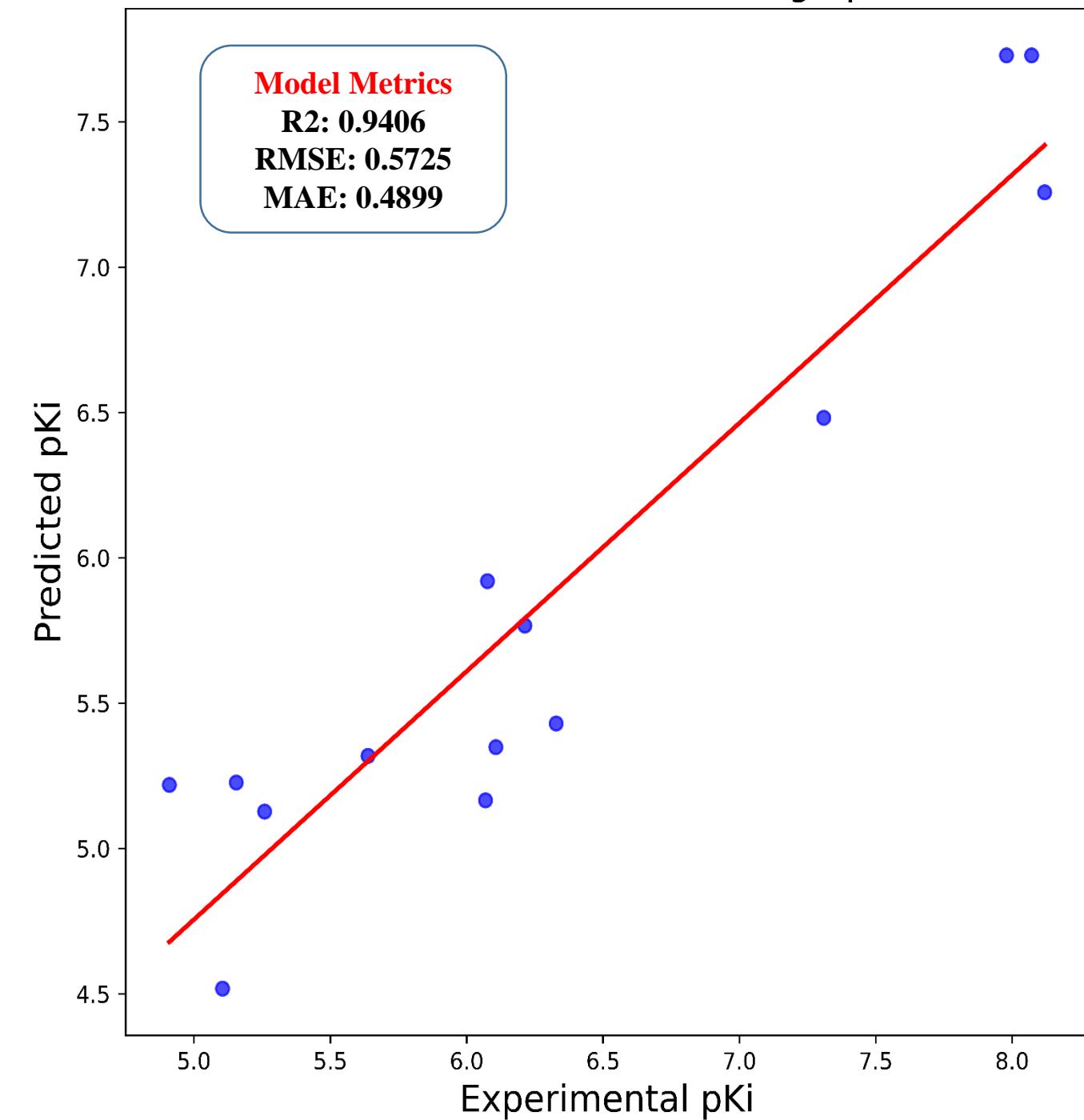
Validation Metrics



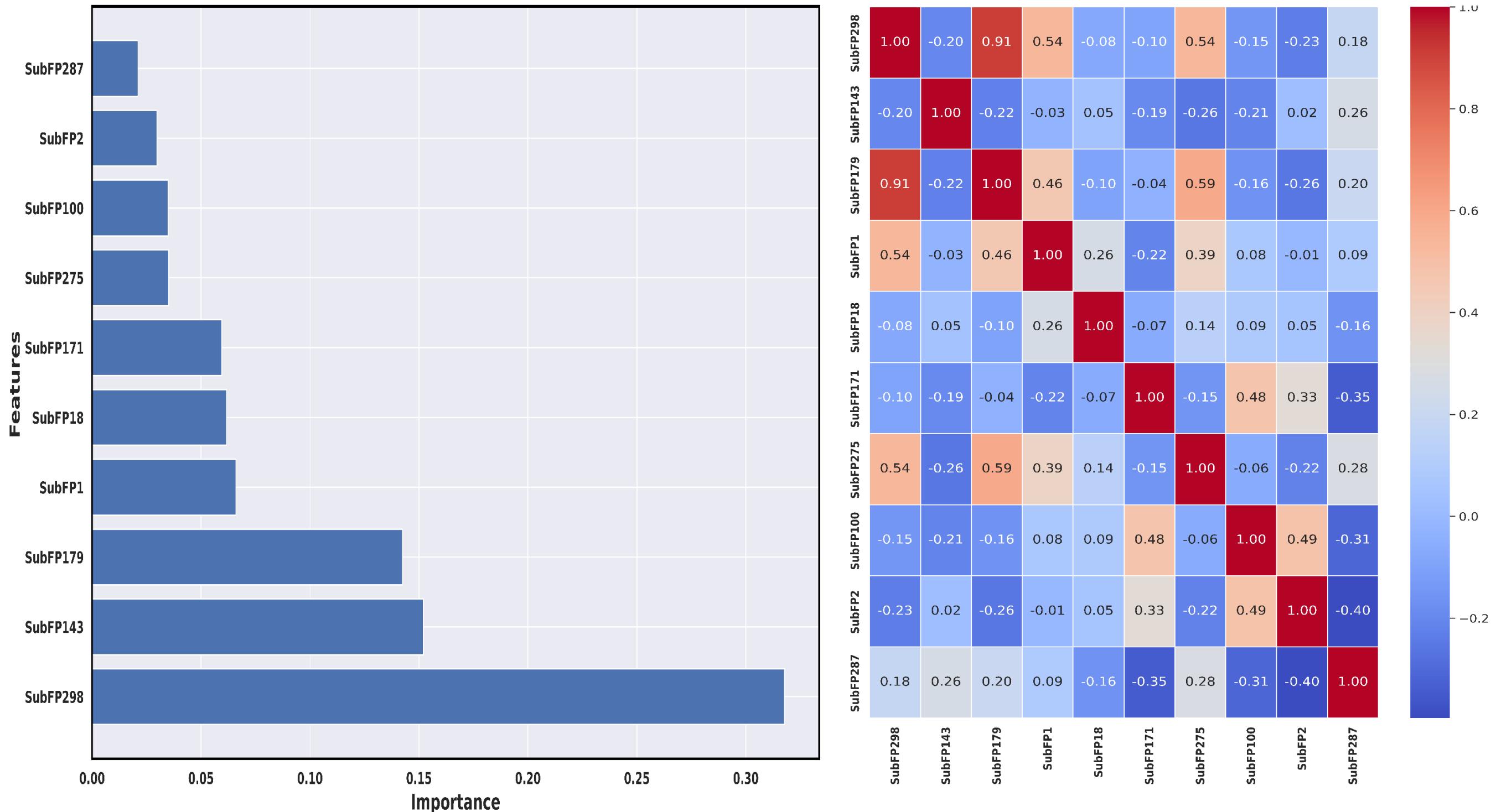
Training set for MtbCA1 substructure fingerprint model



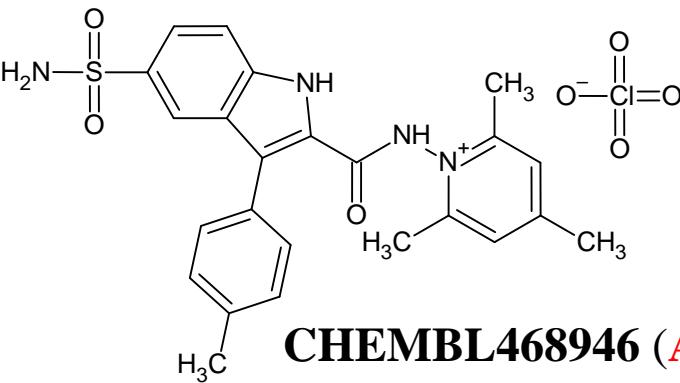
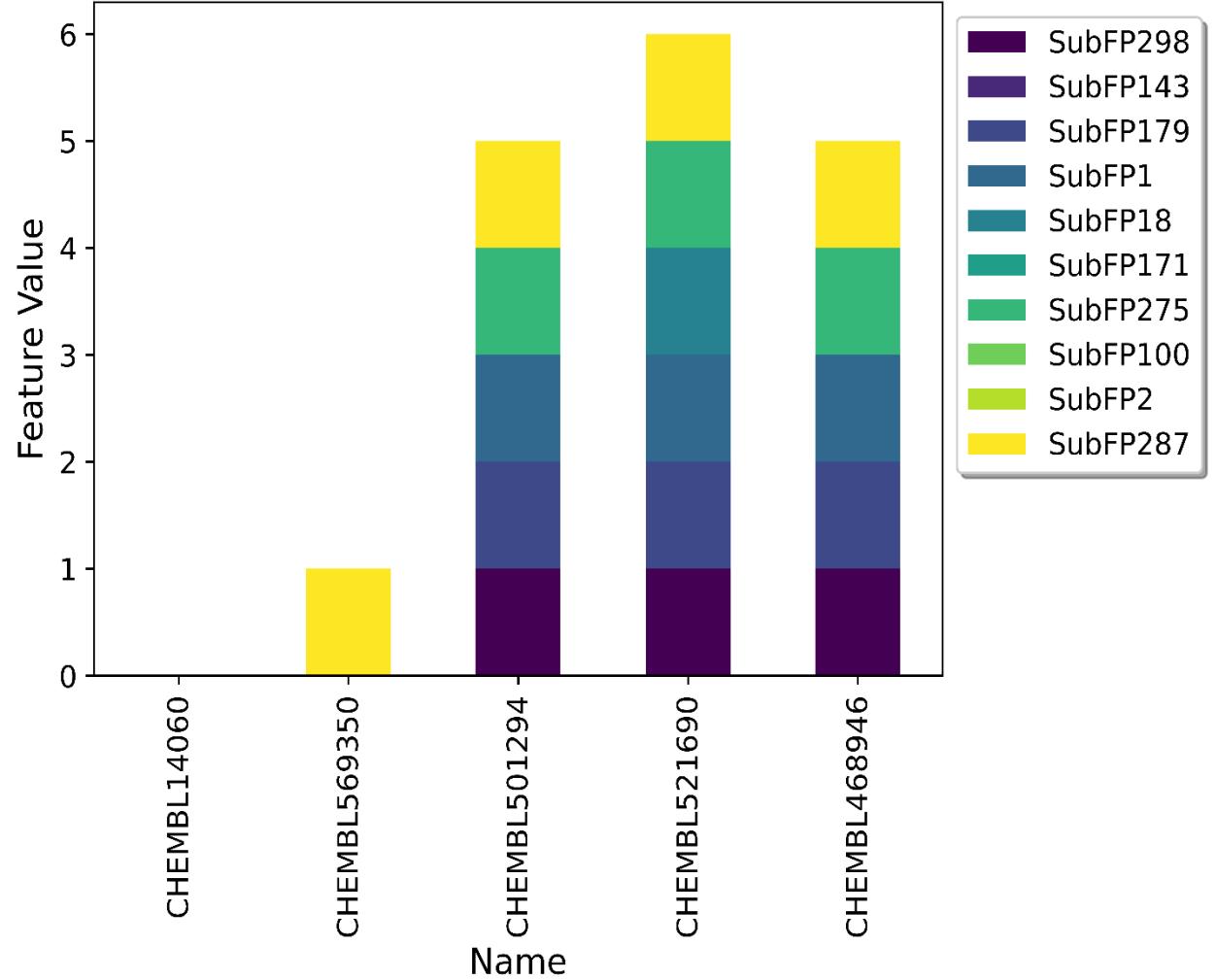
Test set for MtbCA1 substructure fingerprint model



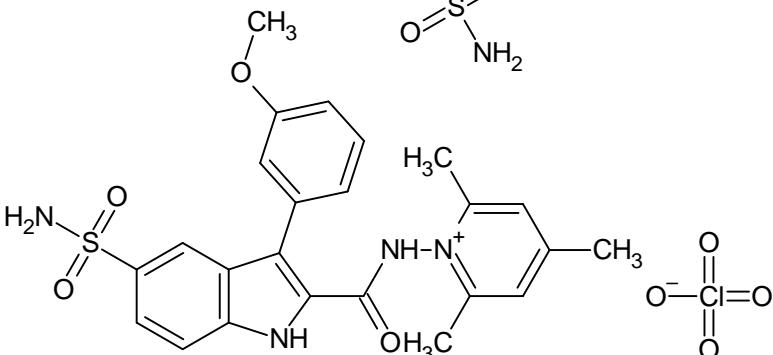
VIP plot and correlation matrix analysis for *MtbCA1* substructure prediction model



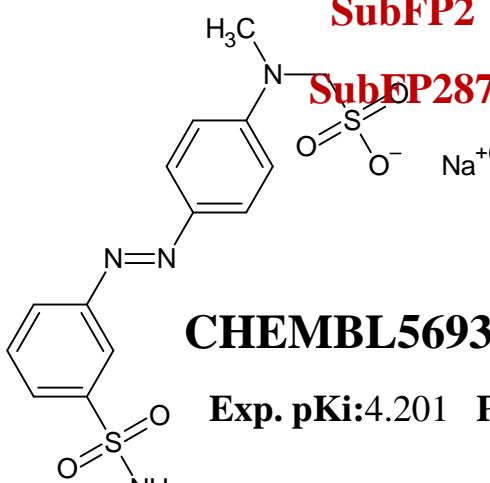
Comparison of substructure fingerprints



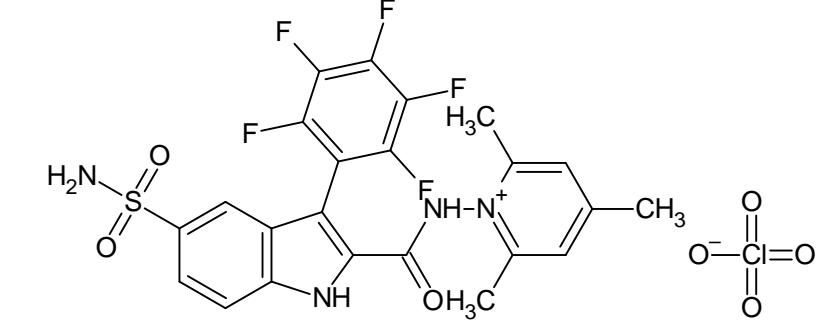
Exp. pKi:8.292 Pred. pKi:7.729



Exp. pKi:8.066 Pred. pKi:7.941



Exp. pKi:4.201 Pred. pKi:4.518



Exp. pKi:4.194 Pred. pKi:4.706

CHEMBL501294 (Active)

Exp. pKi:8.013 Pred. pKi:7.729

SubFP298

Cation

SubFP143

Carbonic acid derivatives

SubFP179

Hetero N basic H

SubFP1

Primary carbon

SubFP18

Alkylarylether

SubFP171

Arylchloride

SubFP275

Heterocyclic

SubFP100

Secondary amide

SubFP2

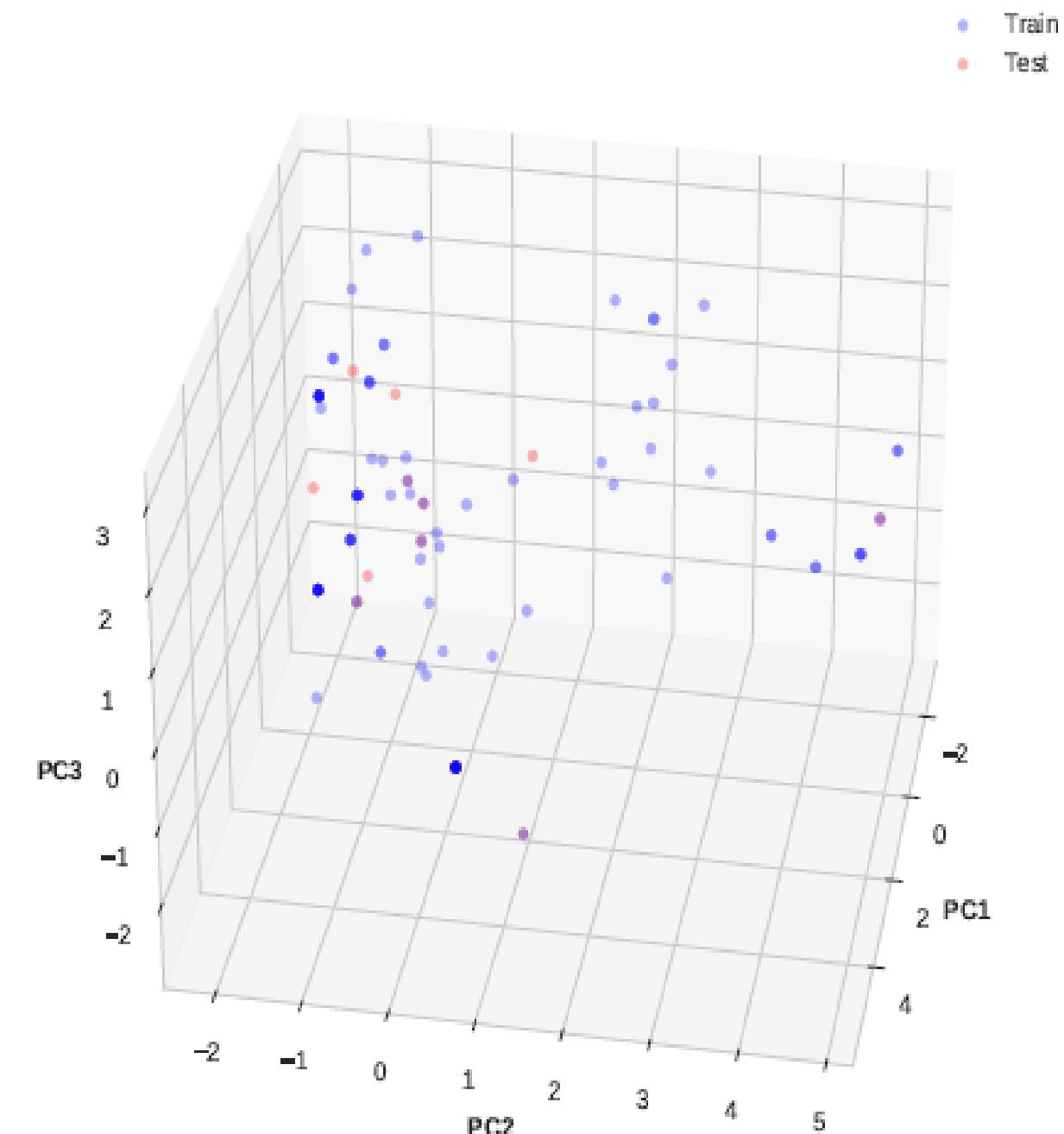
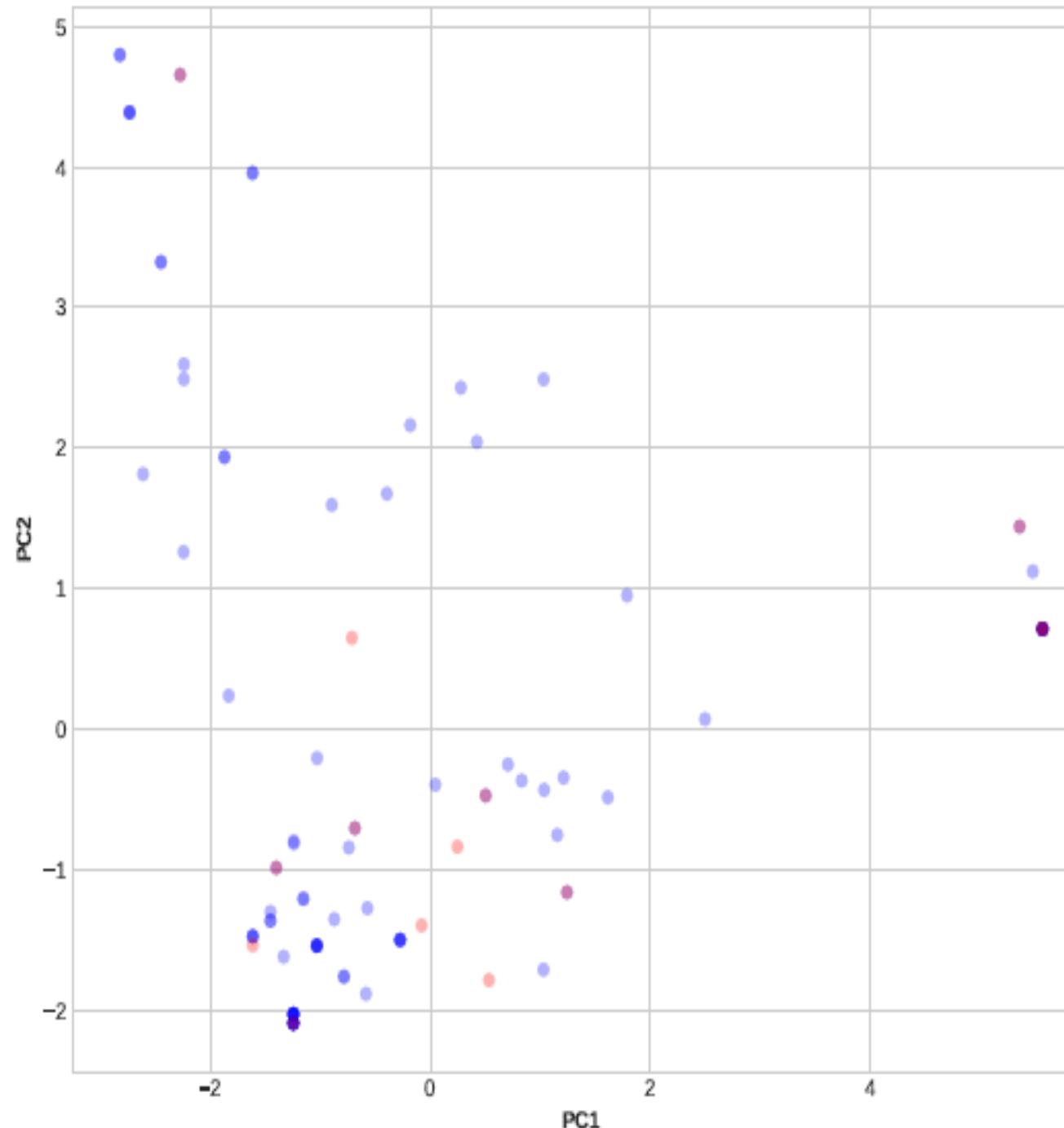
Secondary carbon

SubFP287



CHEMBL14060 (Inactive)

Applicability domain analysis through 2D and 3D PCA plots for *MtbCA1* substructure prediction model



ML-QSAR model for bioactivity prediction of *MtbCA1* inhibitors using 1D and 2D molecular descriptors



Curation of *MtbCA1* inhibitors with Ki value

Total inhibitors: 124

1D and 2D Molecular Descriptors: 1444

Prediction Parameter: Bioactivity Class ('pKi')

Dataset Splitting into training and test set

Elimination of highly correlated and constant fingerprints

Final number of descriptors(763), number of molecules in training(99) and test(25)

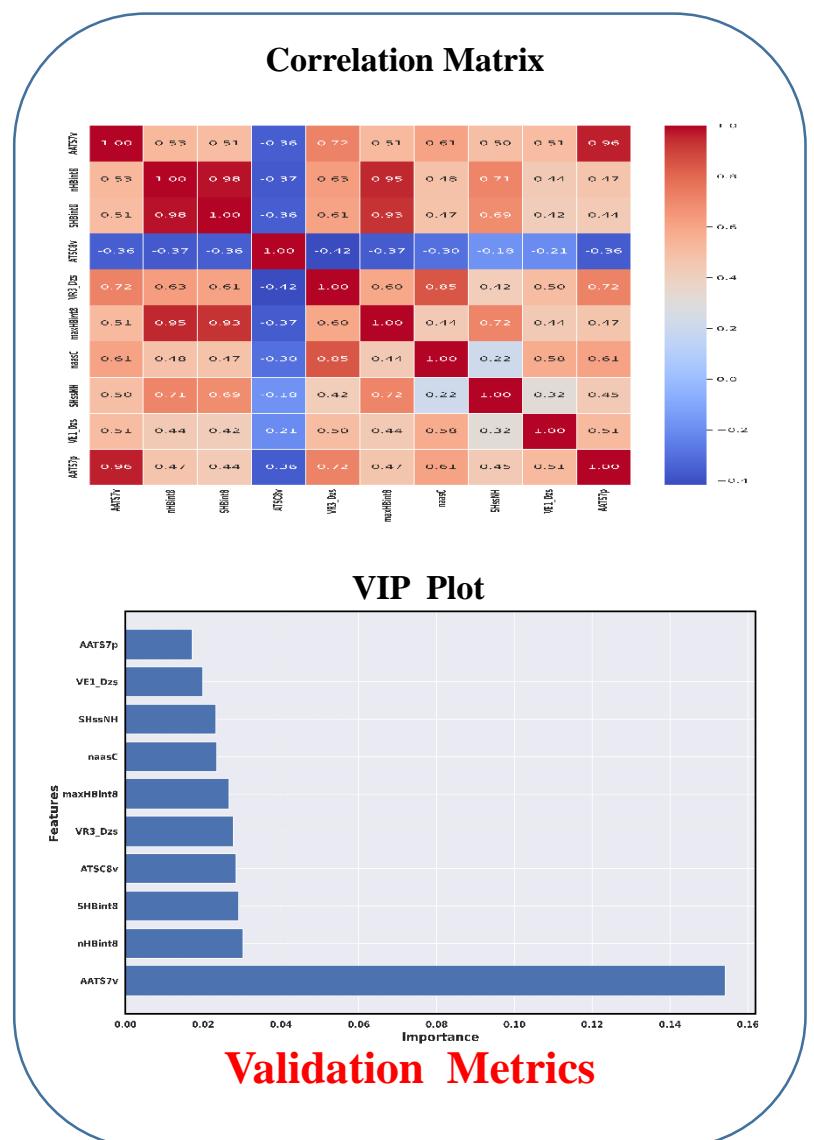
Random Forest Regressor

Initial Model Performance

Model Metrics
R2 (Train: 0.9769, Test: 0.6047)
RMSE (Train: 0.2793, Test: 0.9688)
MAE (Train: 0.2134, Test: 0.7608)

FP and FN Removal

Mention molecules removed (9)

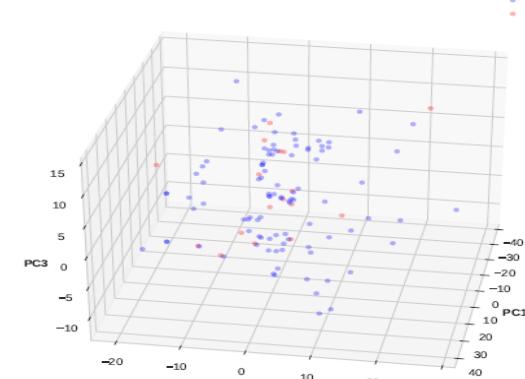


Model Metrics
R2 (Train: 0.9787, Test: 0.9333)
RMSE (Train: 0.2551, Test: 0.4083)
MAE (Train: 0.2042, Test: 0.4899)

Final Model Performance

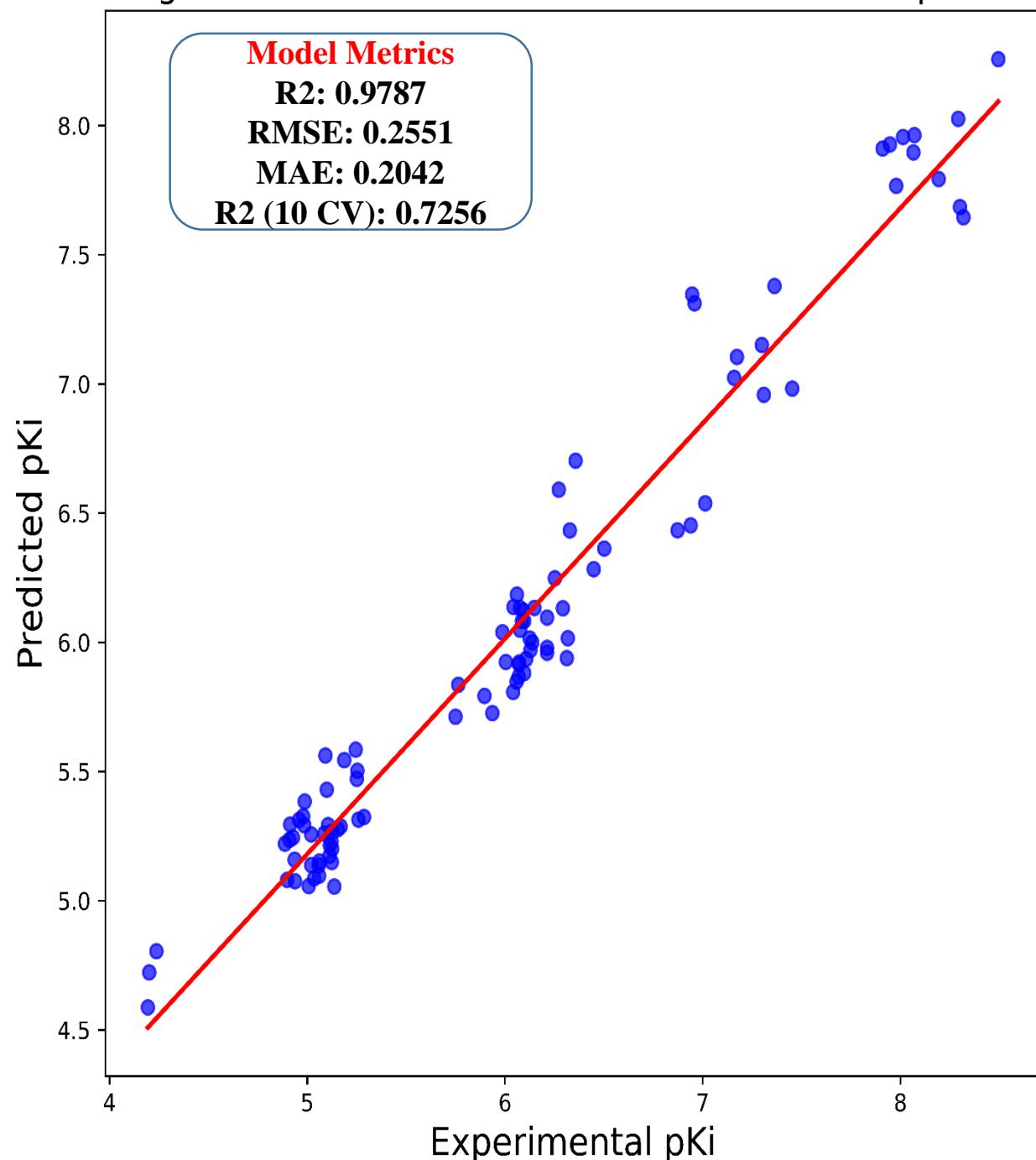
The final machine-learning model

number of molecules in training(98) and test(17)

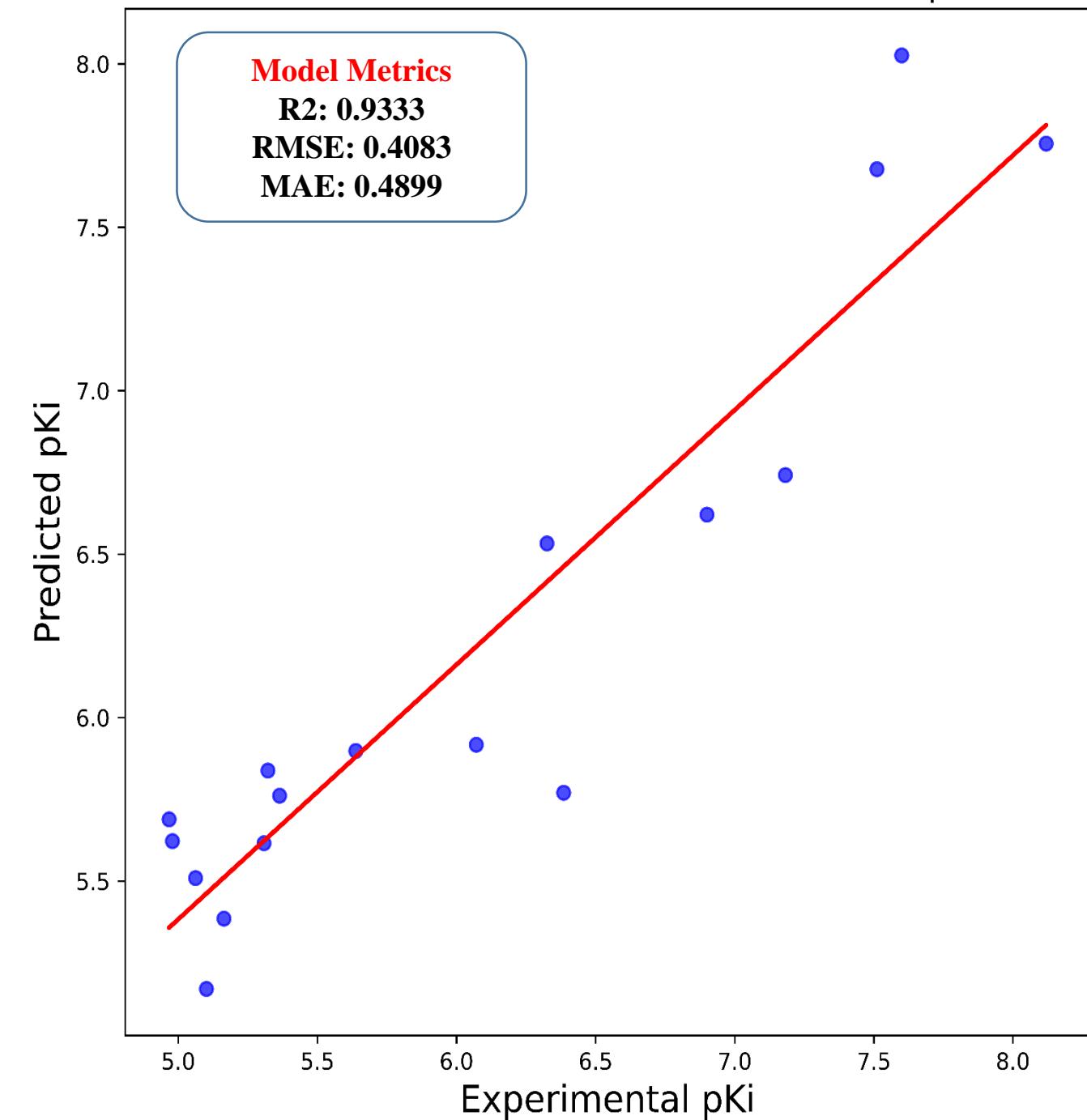


Applicability Domain Analysis

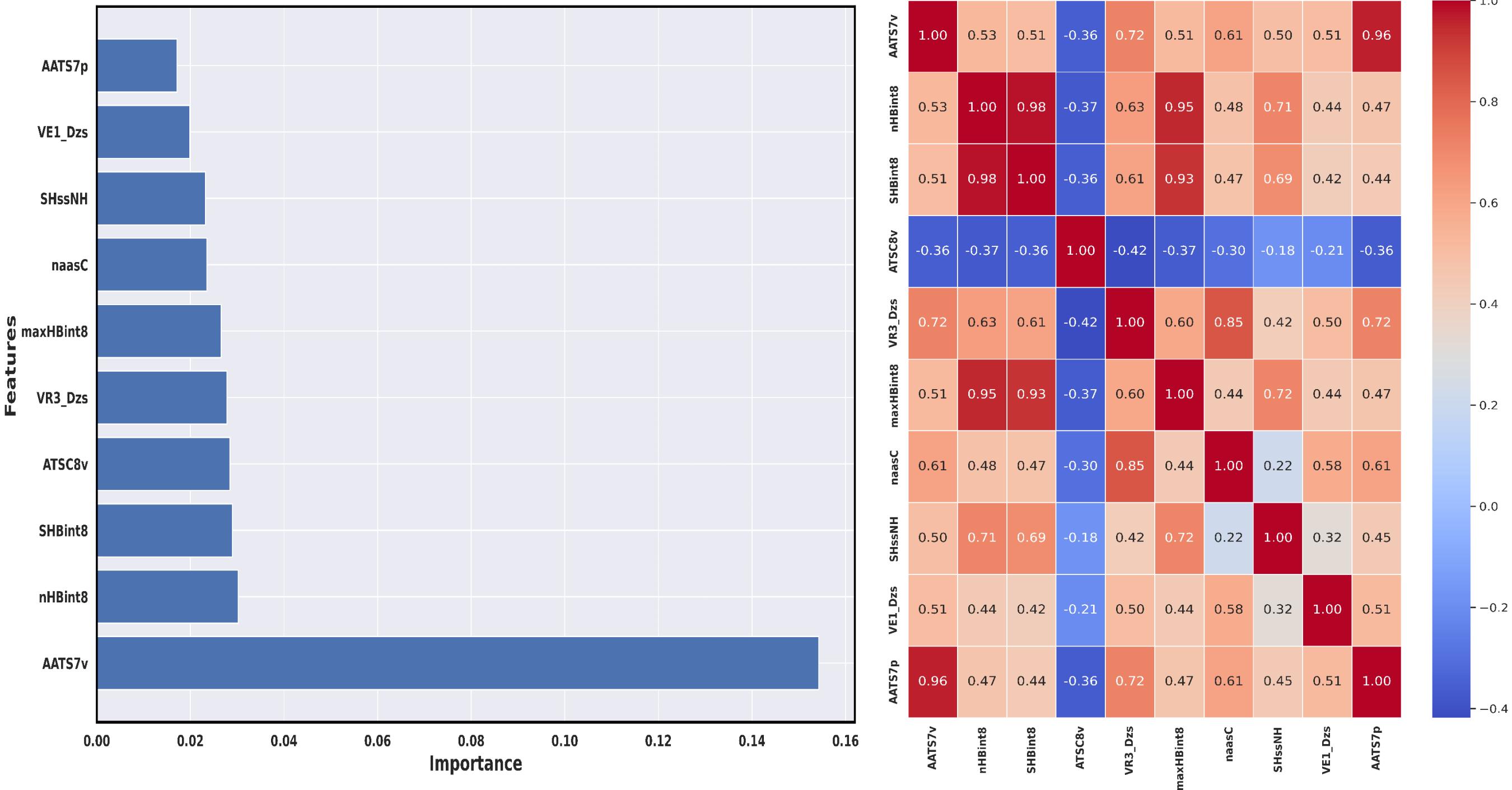
Training set for MtbCA1 1D and 2D molecular descriptor model



Test set for MtbCA1 1D and 2D molecular descriptor model

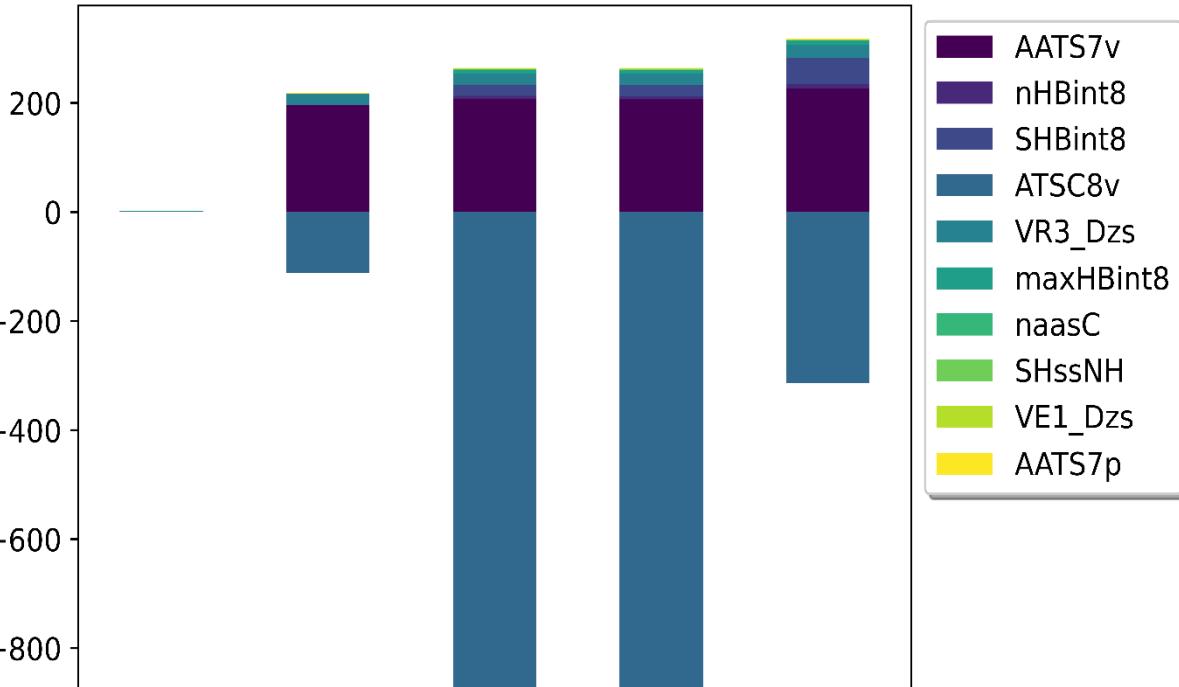


VIP plot and correlation matrix analysis for *MtbCA1* 1D 2D molecular descriptor prediction model



Comparison of 1D 2D molecular descriptors

Feature Value



AATS7v	Average Broto-Moreau autocorrelation - lag 7 / weighted by van der Waals volumes
nHBint8	Count of E-State descriptors of strength for potential Hydrogen Bonds of path length 8
SHBint8	Sum of E-State descriptors of strength for potential hydrogen bonds of path length 8
ATSC8v	Centered Broto-Moreau autocorrelation - lag 8 / weighted by van der Waals volumes
VR3_Dzs	Logarithmic Randic-like eigenvector-based index from Barysz matrix / weighted by I-state
maxHBint8	Maximum E-State descriptors of strength for potential Hydrogen Bonds of path length 8
naasC	Count of atom-type E-State: :C:-
SHssNH	Sum of atom-type H E-State: -NH-
VE1_Dzs	Coefficient sum of the last eigenvector from Barysz matrix / weighted by I-state
AATS7p	Average Broto-Moreau autocorrelation - lag 7 / weighted by polarizabilities

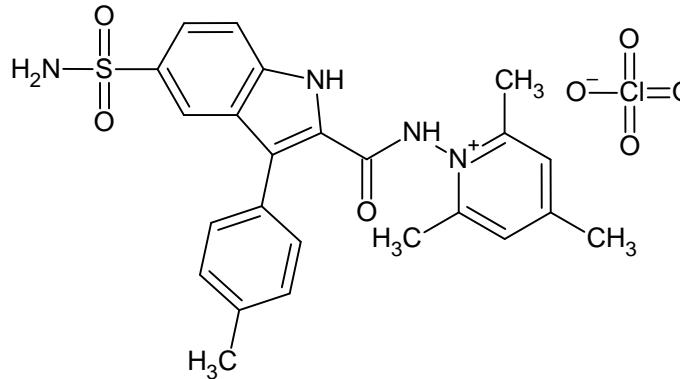
CHEMBL14060

CHEMBL569350

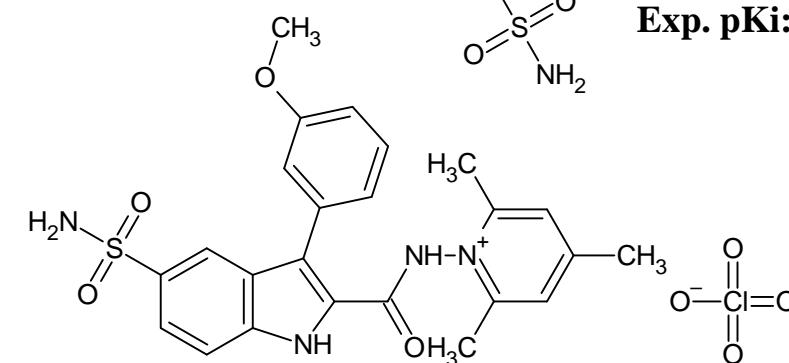
CHEMBL468946

CHEMBL521690

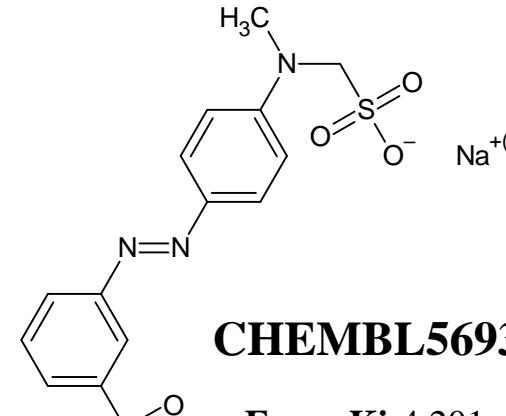
CHEMBL501294



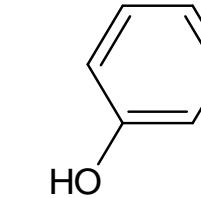
Exp. pKi:8.292 Pred. pKi:8.026



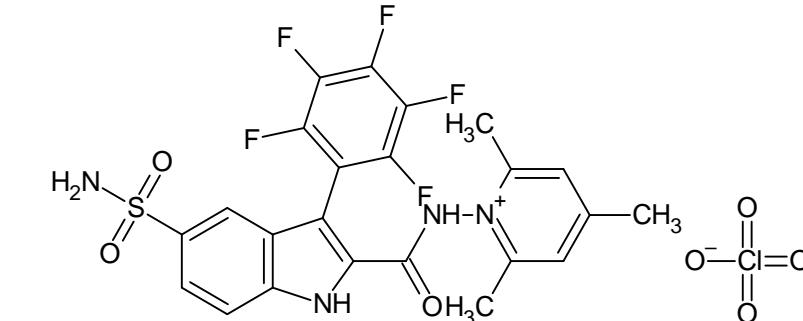
Exp. pKi:8.066 Pred. pKi:7.896



Exp. pKi:4.201 Pred. pKi:4.723

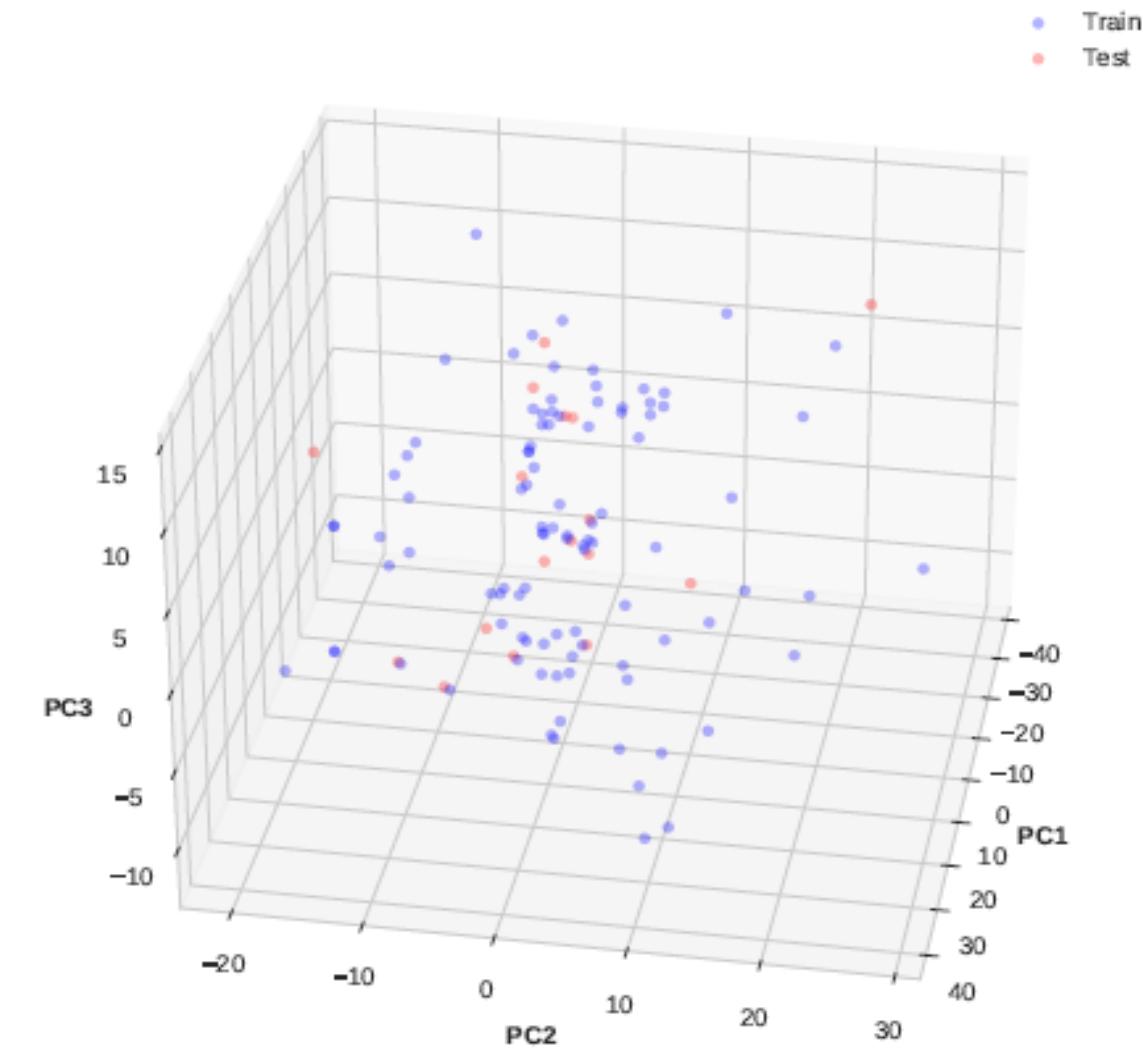
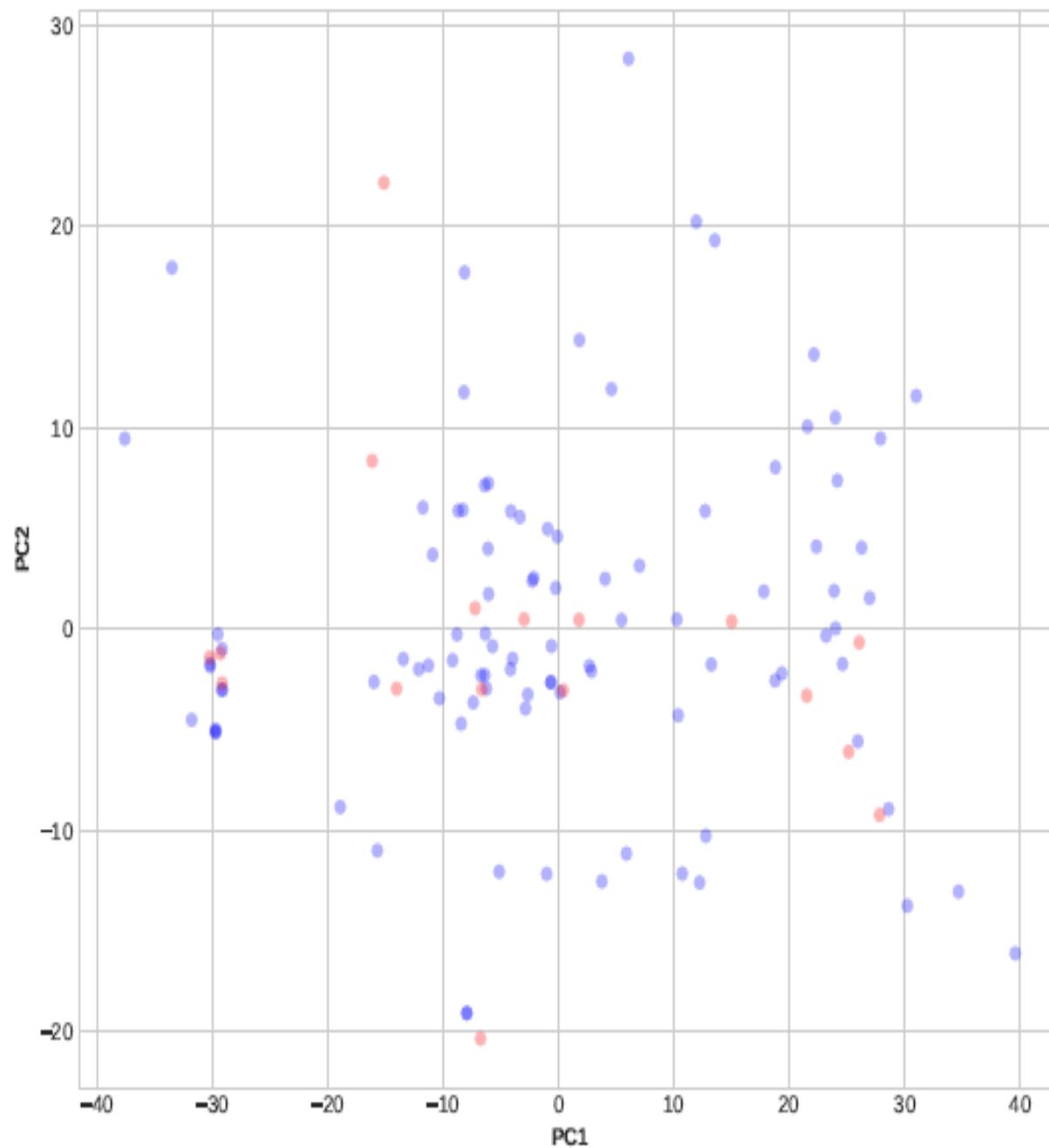


Exp. pKi:4.194 Pred. pKi:4.587



Exp. pKi:8.013 Pred. pKi:7.956

Applicability domain analysis through 2D and 3D PCA plots for *MtbCA1* 1D 2D molecular descriptor prediction model



MtbCA2 Prediction Model

-5.000

5.000

ML-QSAR model for bioactivity prediction of *MtbCA2* inhibitors using substructure fingerprints



Curation of *MtbCA2* inhibitors with Ki value

Total inhibitors: 42

Substructure fingerprints: 307

Prediction Parameter: Bioactivity Class ('pKi')

Dataset Splitting into training and test set

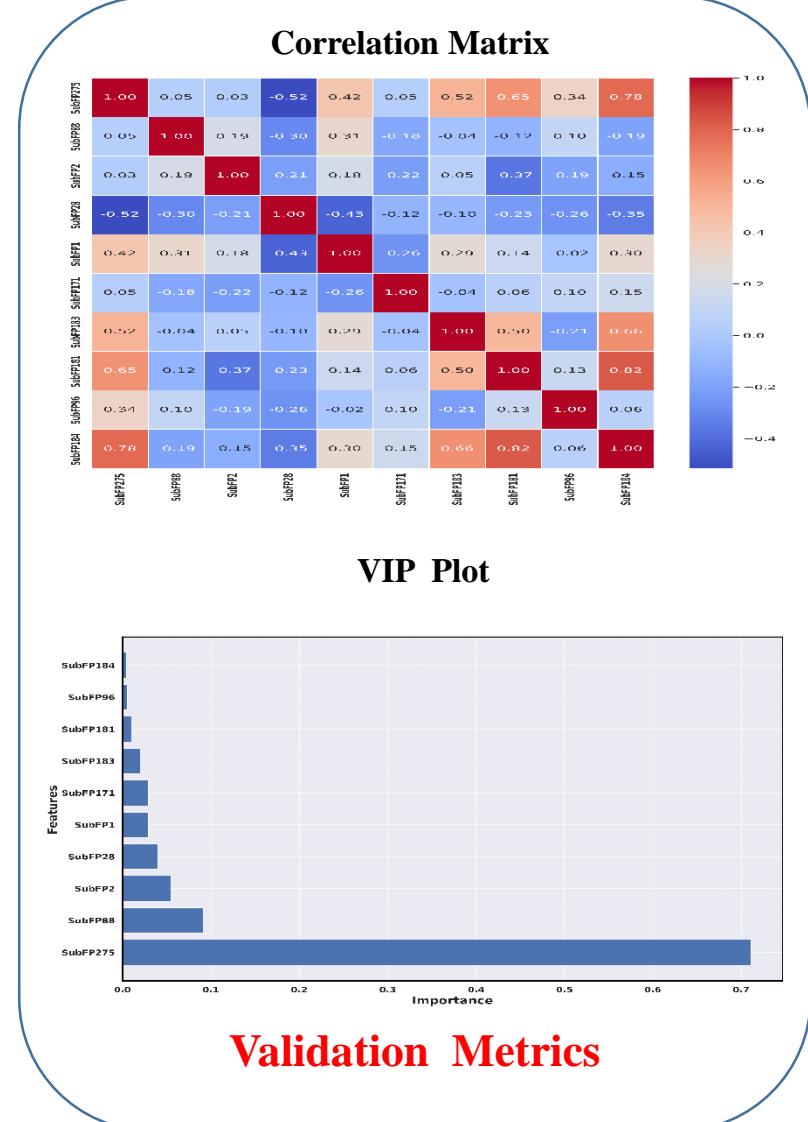
Elimination of highly correlated and constant molecular descriptors

Final number of fingerprints(13), number of molecules in training(33) and test(9)

Random Forest Regressor

Initial Model Performance

Model Metrics
R2 (Train: 0.9645, Test: 0.8764)
RMSE (Train: 0.2785, Test: 0.5463)
MAE (Train: 0.1948, Test: 0.4097)

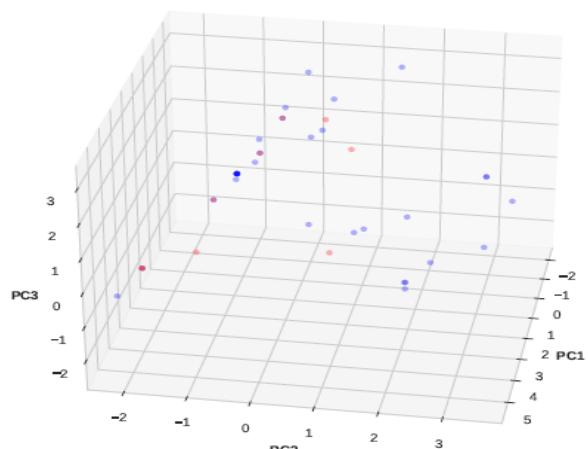


Model Metrics
R2 (Train: 0.9645, Test: 0.8764)
RMSE (Train: 0.2785, Test: 0.5463)
MAE (Train: 0.1948, Test: 0.4097)

Final Model Performance

The final machine-learning model

number of molecules in training(33) and test(9)

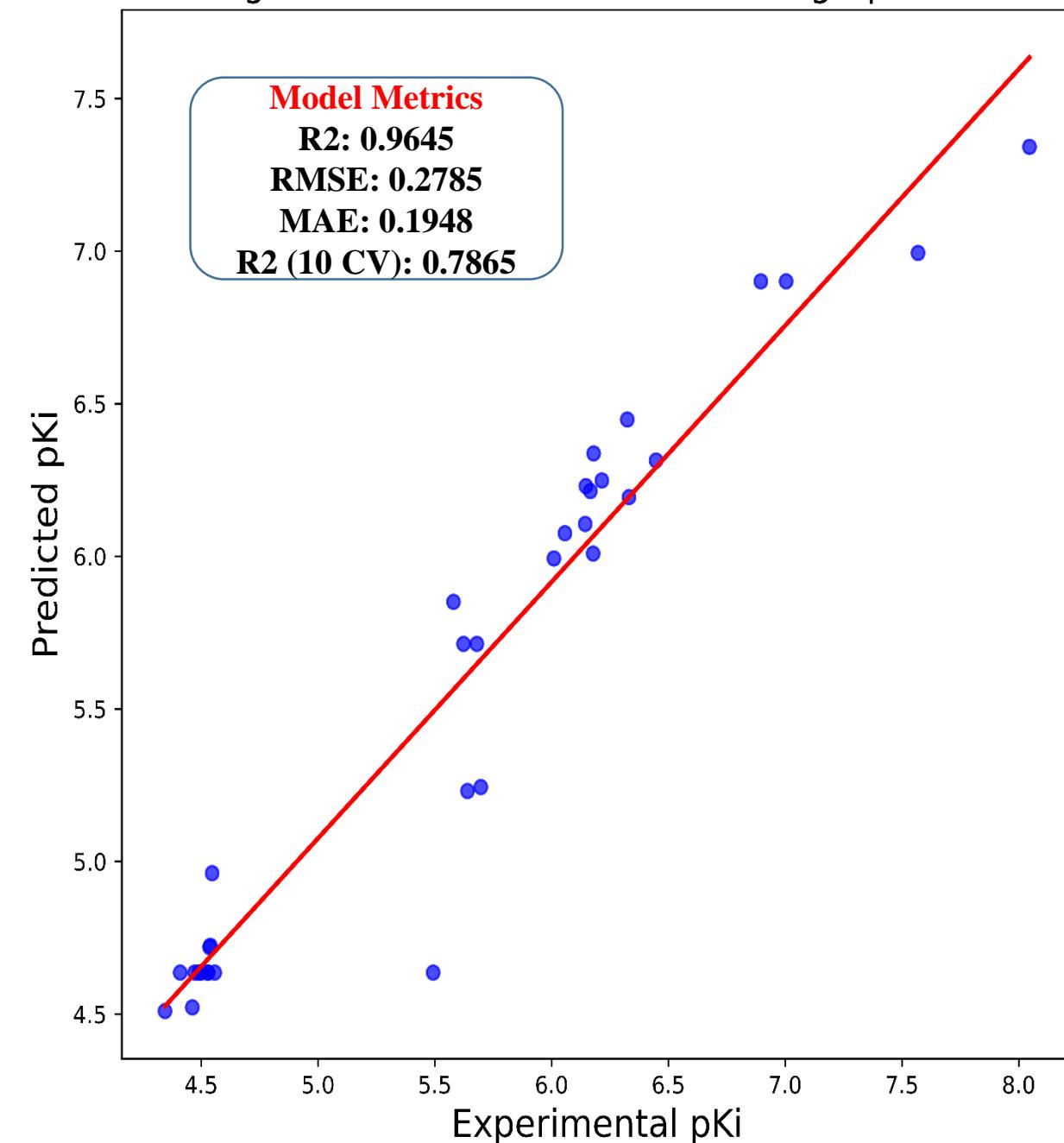


FP and FN Removal

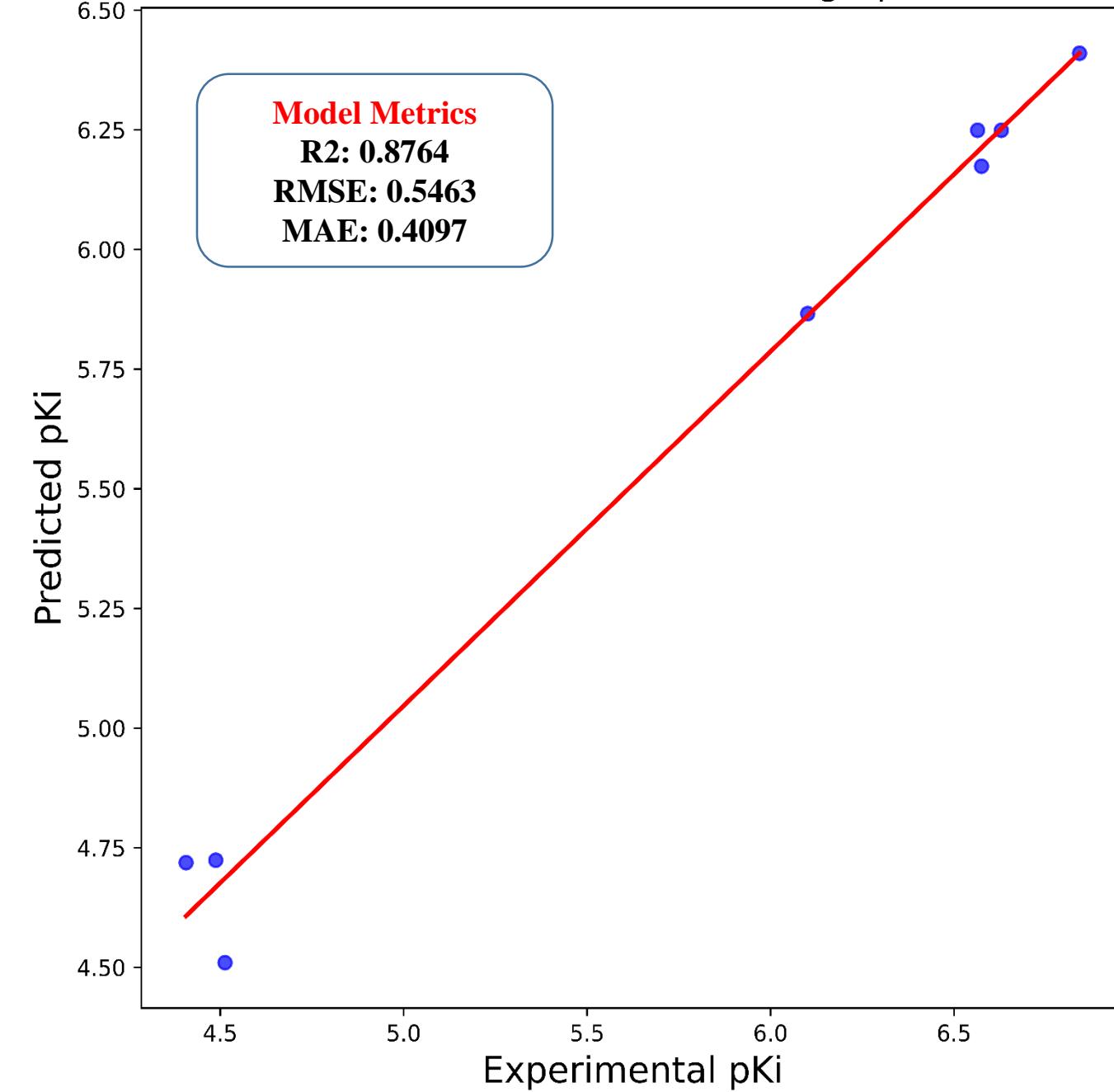
Mention molecules removed (0)

Applicability Domain Analysis

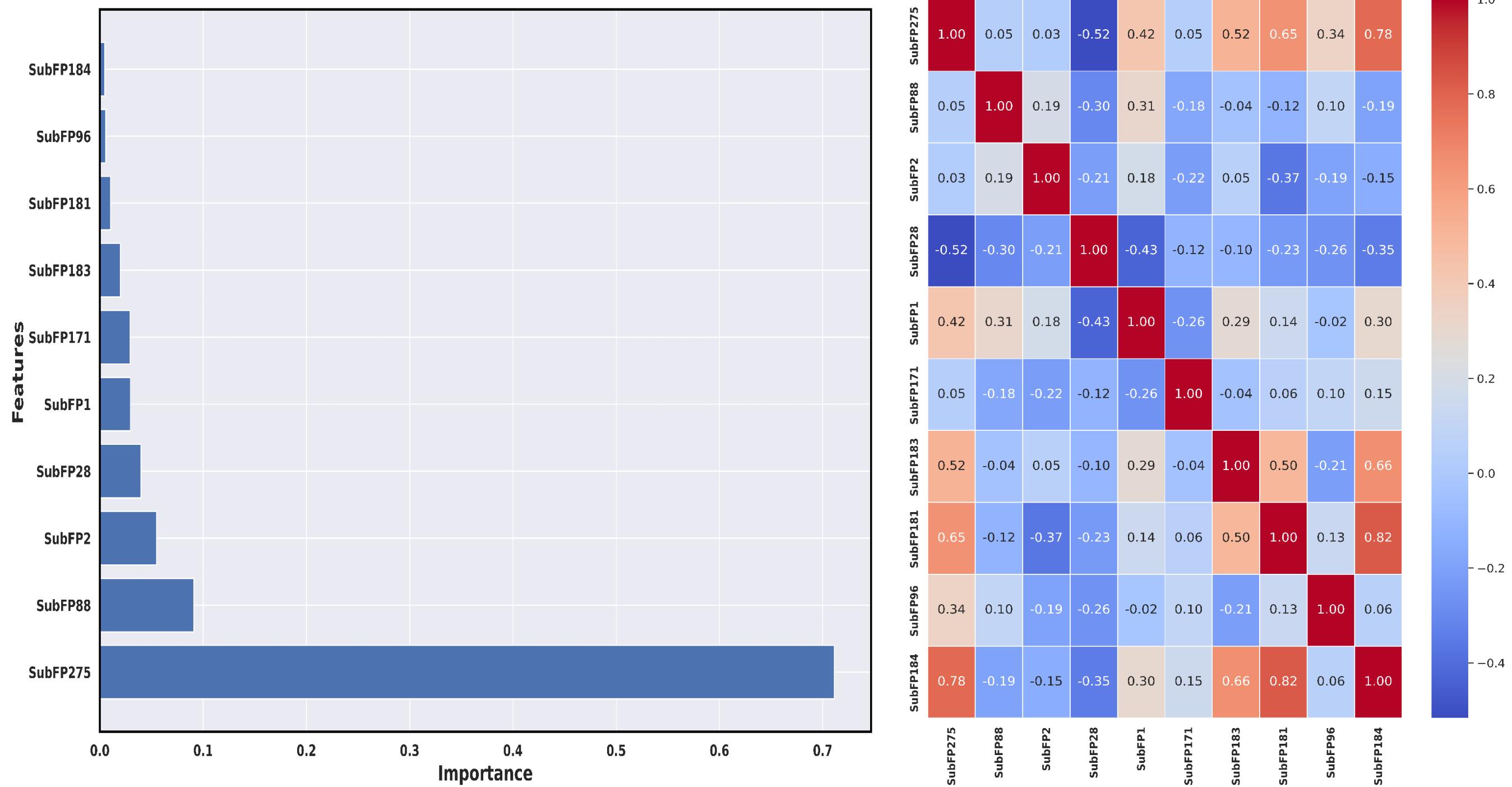
Training set for MtbCA2 substructure fingerprint model



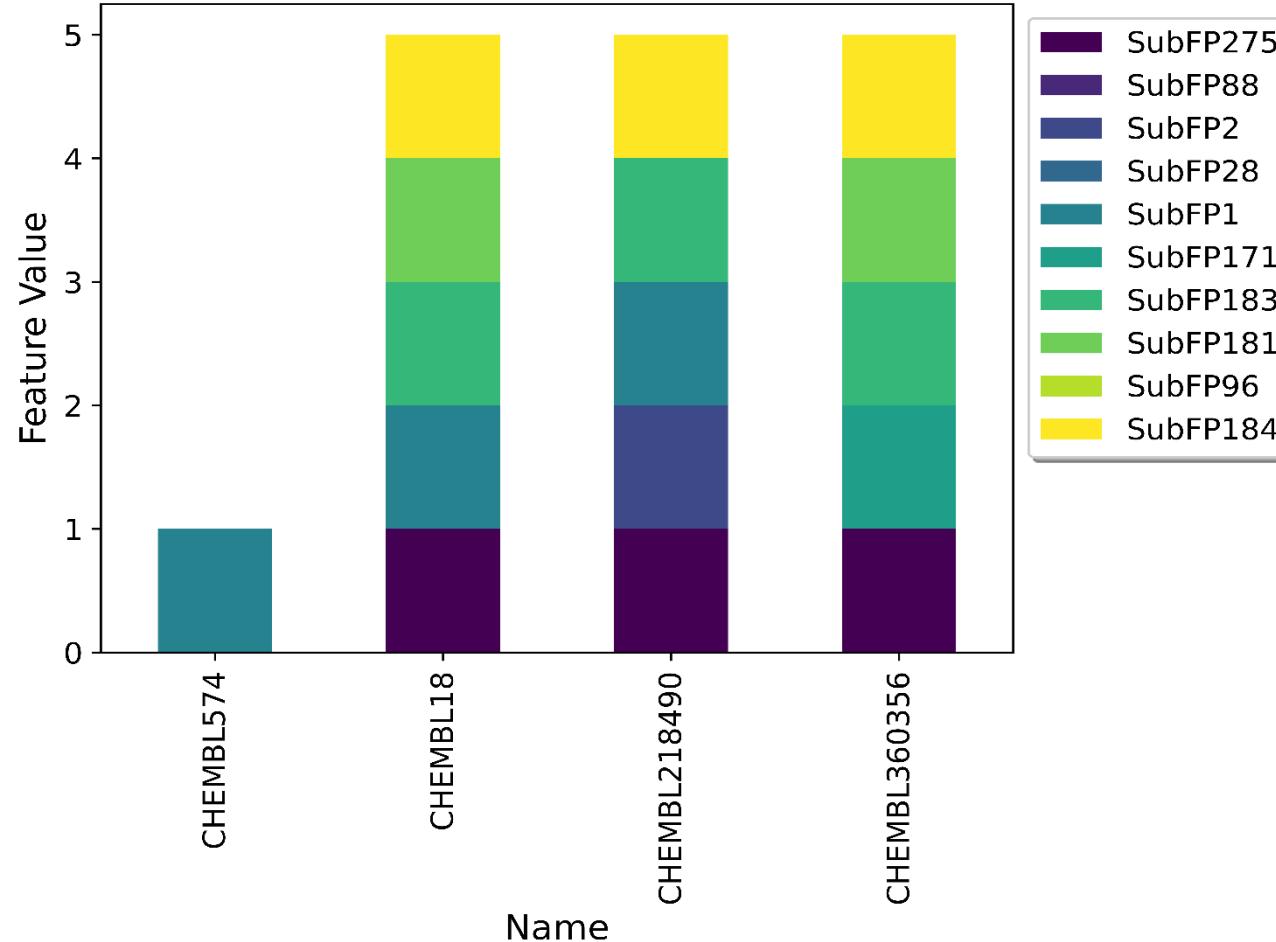
Test set for MtbCA2 substructure fingerprint model



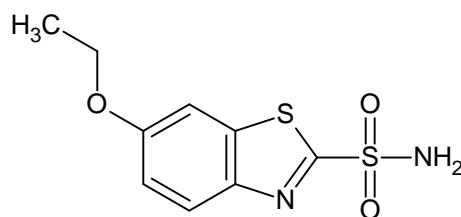
VIP plot and correlation matrix analysis for *MtbCA2* substructure prediction model



Comparison of substructure fingerprints

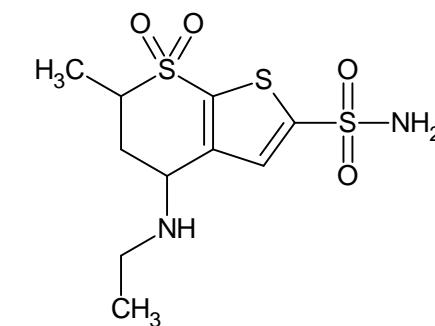


SubFP275	Heterocyclic
SubFP88	Carboxylic acid derivative
SubFP2	Secondary carbon
SubFP28	Primary aromatic amine
SubFP1	Primary carbon
SubFP171	Arylchloride
SubFP183	Hetero S
SubFP181	Hetero N nonbasic
SubFP96	Carbodithioic ester
SubFP184	Heteroaromatic



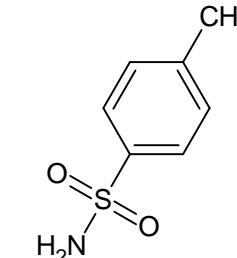
**CHEMBL18/Ethoxzolamide
(Active)**

Exp. pKi:7.569 Pred. pKi:6.994



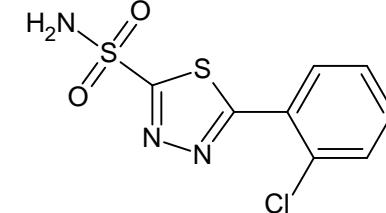
**CHEMBL218490/Dorzolamide
(Active)**

Exp. pKi:7.004 Pred. pKi:6.901



**CHEMBL574/P-Toluenesulfonamide
(Inactive)**

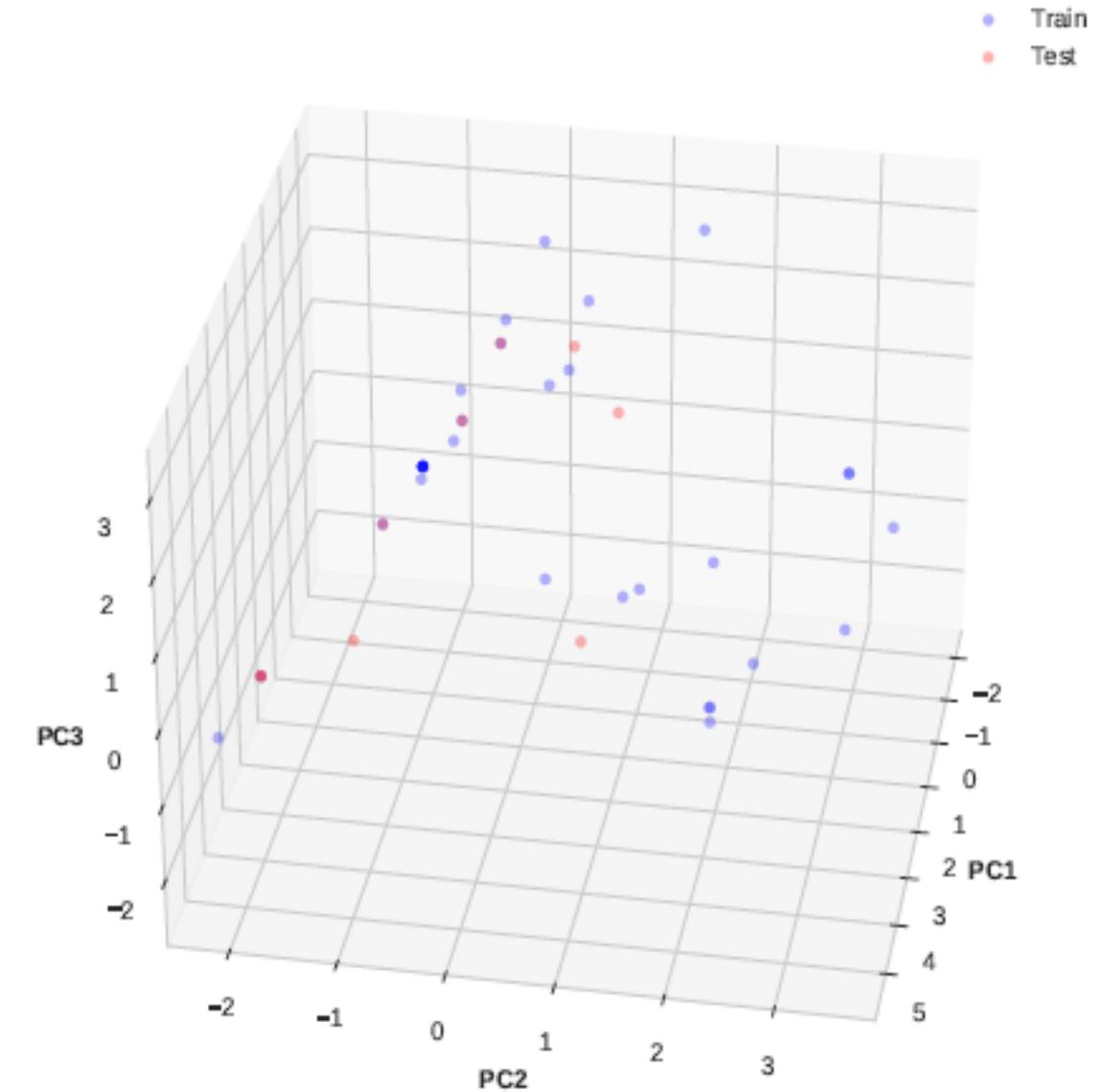
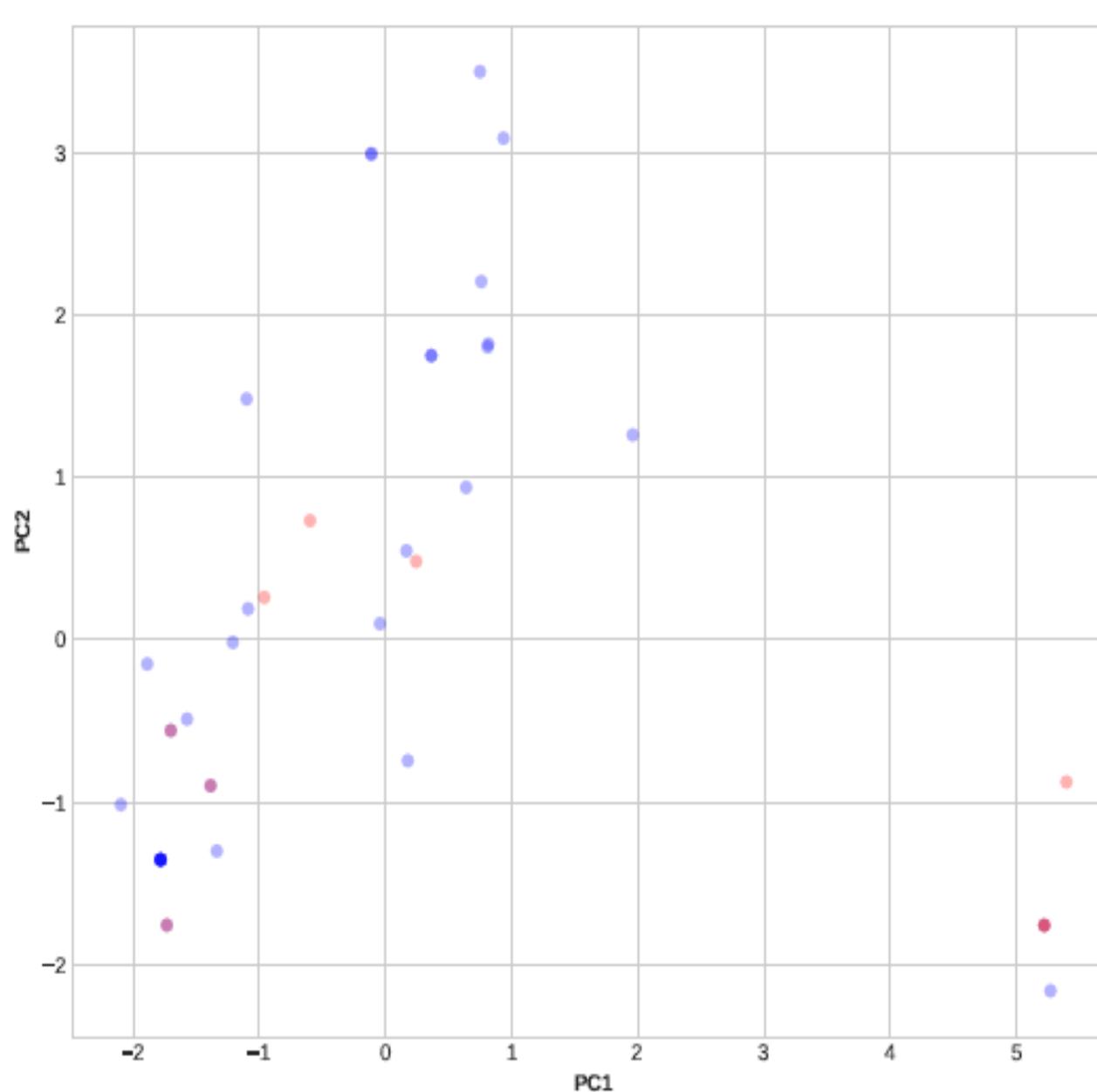
Exp. pKi:4.41 Pred. pKi:4.636



**CHEMBL360356/Chlorzolamide
(Inactive)**

Exp. pKi:4.345 Pred. pKi:4.51

Applicability domain analysis through 2D and 3D PCA plots for *MtbCA2* substructure prediction model



ML-QSAR model for bioactivity prediction of *MtbCA2* inhibitors using 1D and 2D molecular descriptors



Curation of *MtbCA2* inhibitors with Ki value

Total
inhibitors:
42

1D and 2D
Molecular
Descriptors:
1444

Prediction
Parameter:
Bioactivity Class
(‘pKi’)

Dataset
Splitting into
training and
test set

Elimination of highly
correlated and constant
molecular descriptors

Final number of
descriptors(736),
number of molecules in
training(33) and test(9)

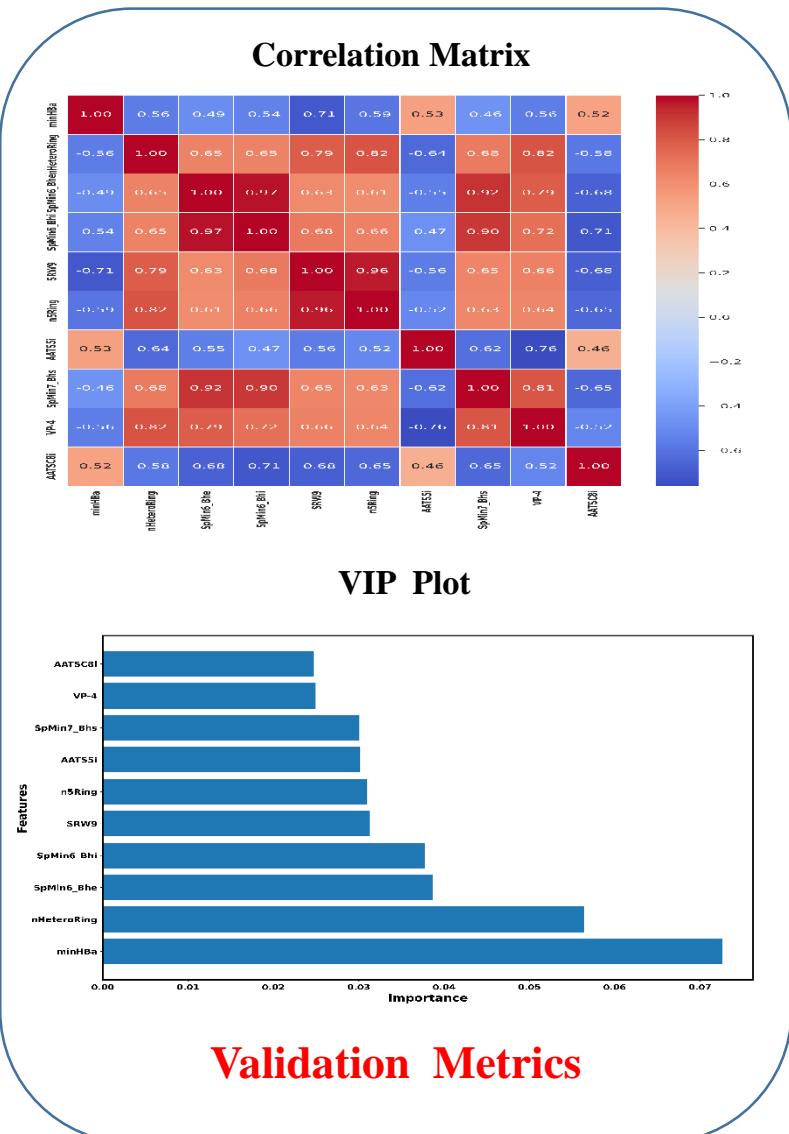
Random Forest Regressor

Initial Model
Performance

Model Metrics
R2 (Train: 0.9651,
Test: 0.8572)
RMSE (Train: 0.326,
Test: 0.4517)
MAE (Train: 0.2082,
Test: 0.3368)

FP and FN Removal

Mention molecules
removed (0)



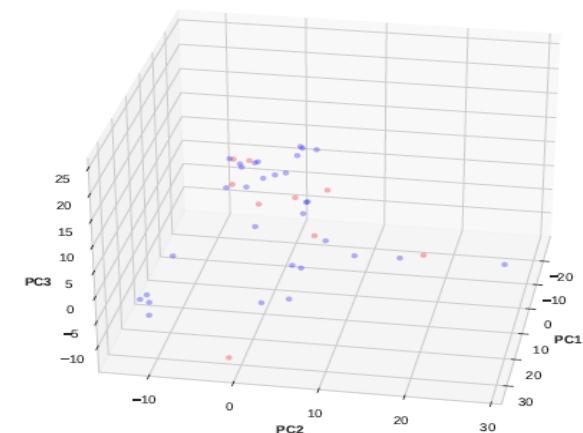
Validation Metrics

Model Metrics
R2 (Train: 0.9864, Test:
0.8874)
RMSE (Train: 0.1971,
Test: 0.4122)
MAE (Train: 0.1553, Test:
0.3101)

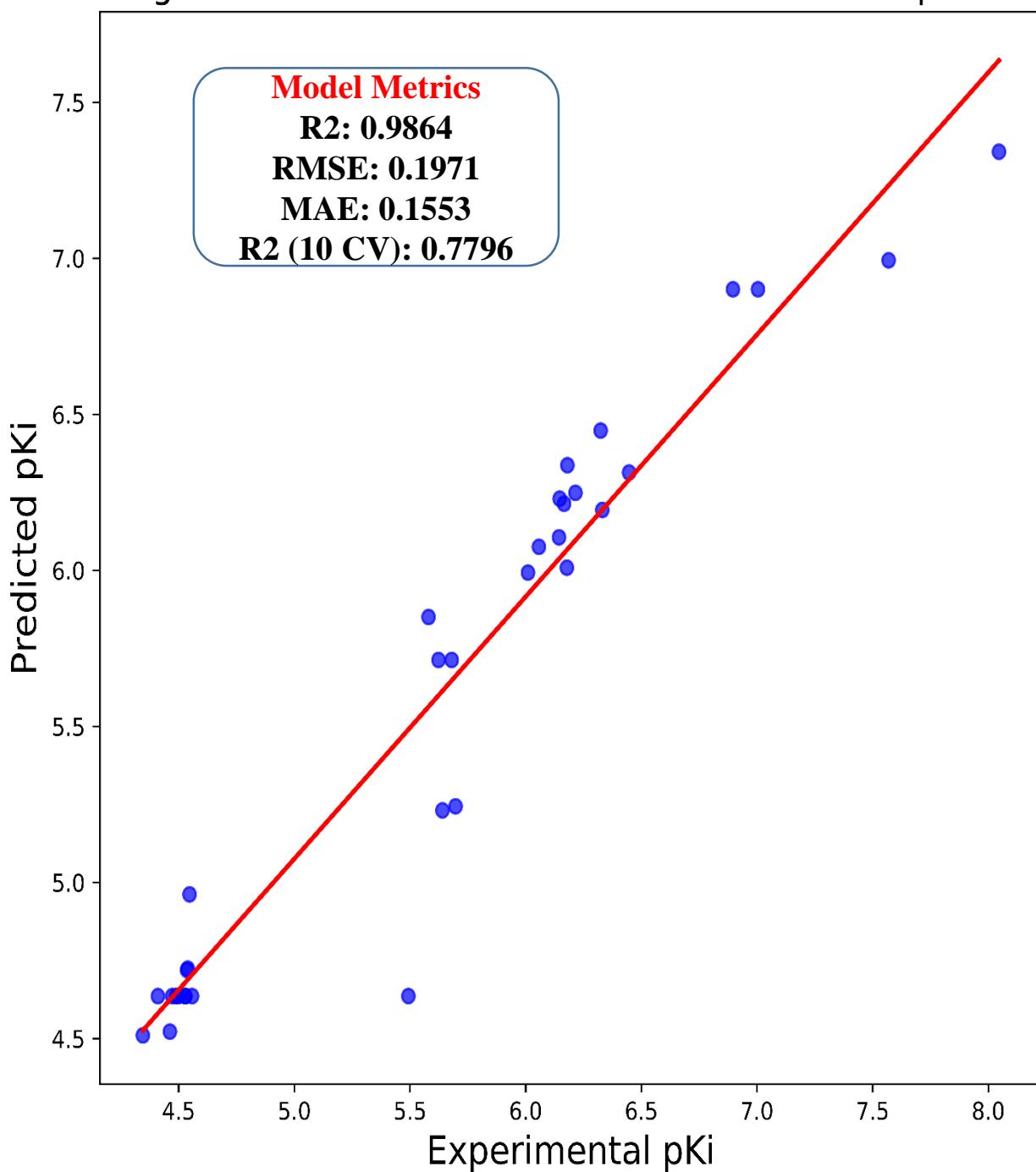
Final Model Performance

The final machine-learning
model

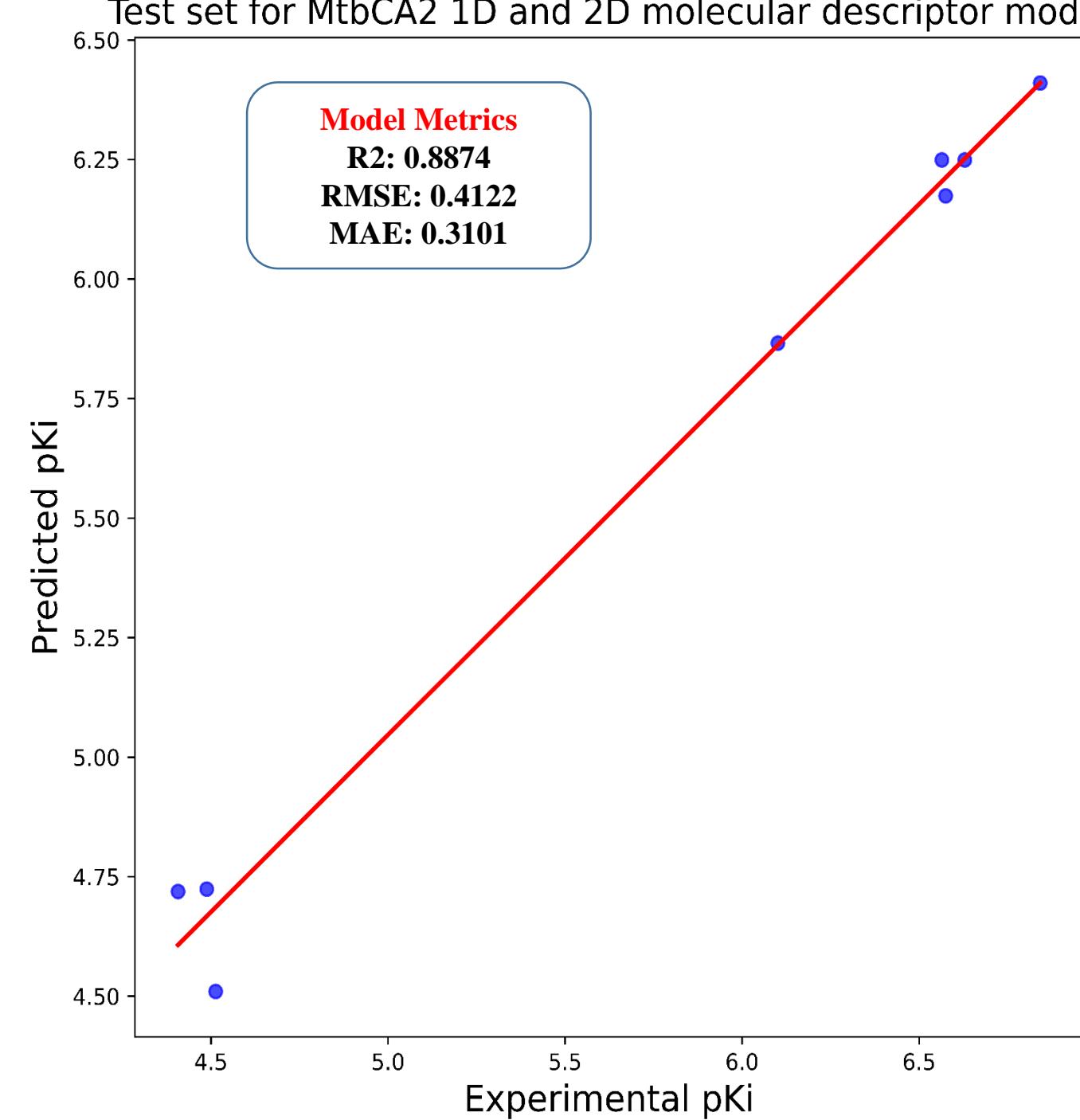
number of
molecules in
training(31) and
test(9)



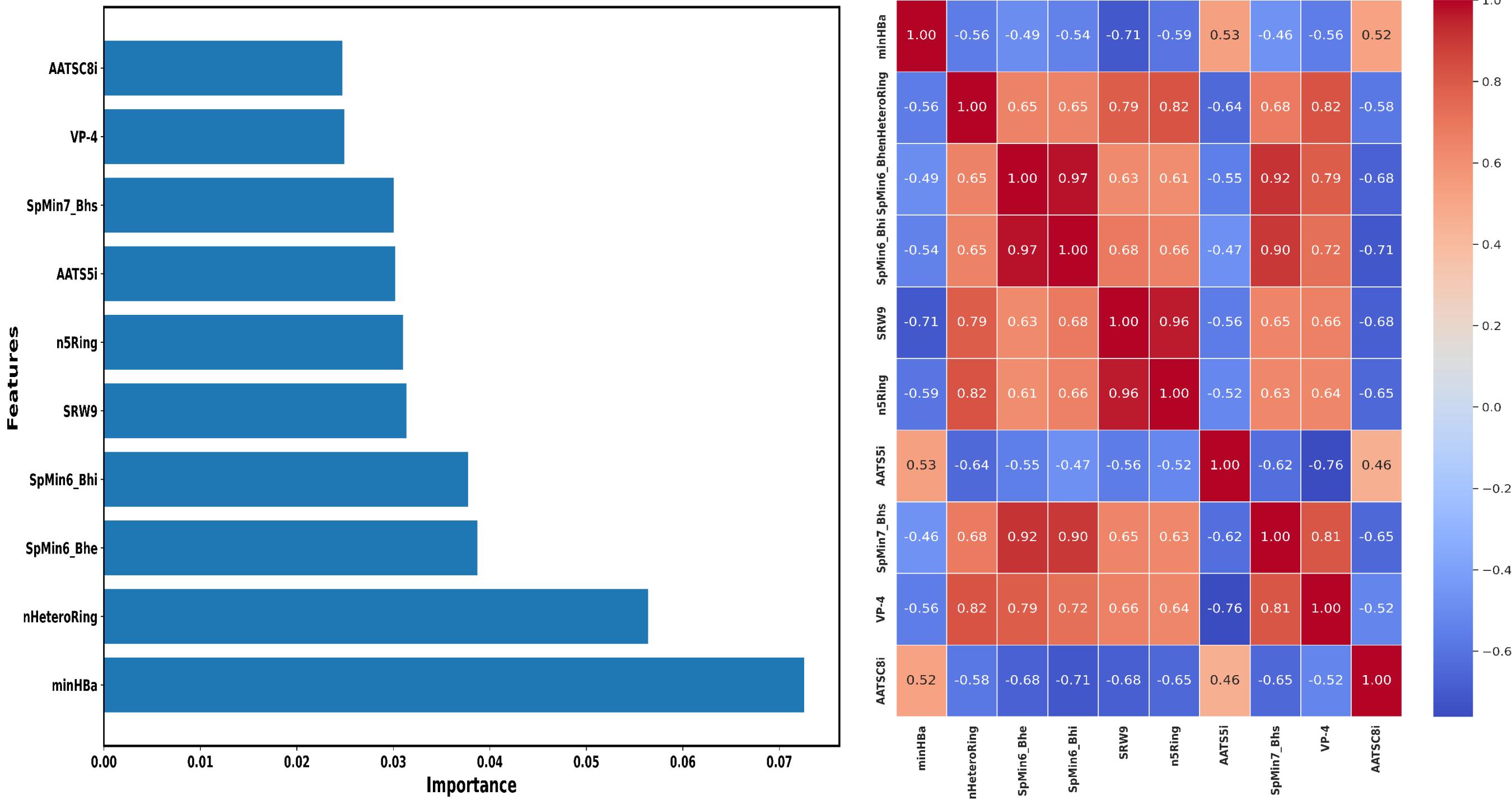
Training set for MtbCA2 1D and 2D molecular descriptor model



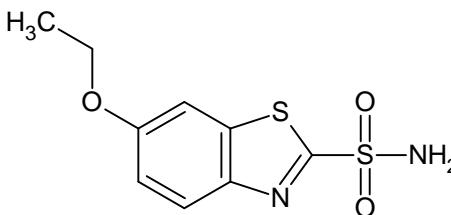
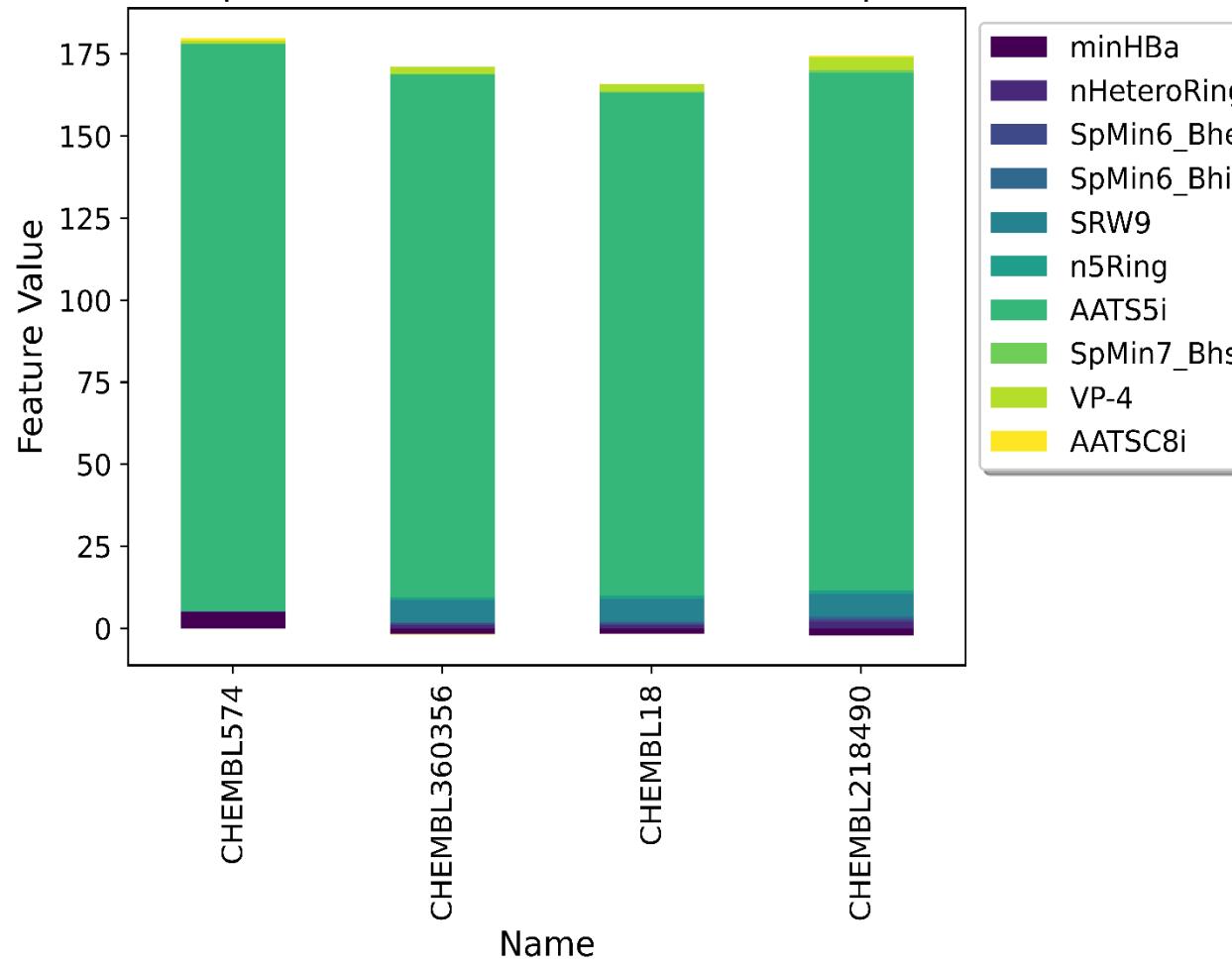
Test set for MtbCA2 1D and 2D molecular descriptor model



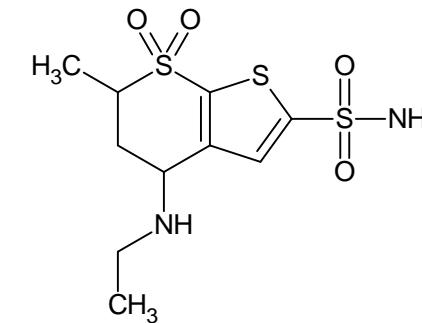
VIP plot and correlation matrix analysis for *MtbCA2* 1D 2D molecular descriptor prediction model



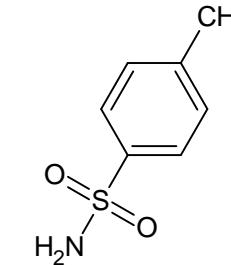
Comparison of 1D 2D molecular descriptors



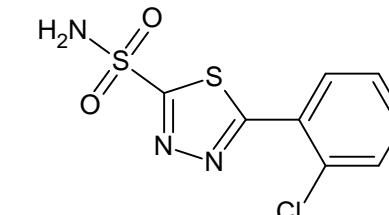
**CHEMBL18/Ethoxzolamide
(Active)**



**CHEMBL218490/Dorzolamide
(Active)**



**CHEMBL574/P-Toluenesulfonamide
(Inactive)**



**CHEMBL360356/Chlorzolamide
(Inactive)**

Exp. pKi:7.569 Pred. pKi:6.823

Exp. pKi:7.004 Pred. pKi:6.784

Exp. pKi:4.41 Pred. pKi:4.907

Exp. pKi:4.345 Pred. pKi:4.491

minHBa

Minimum E-States for (strong) Hydrogen Bond acceptors

nHeteroRing

Number of rings containing heteroatoms (N, O, P, S, or halogens)

SpMin6_Bhe

Smallest absolute eigenvalue of Burden modified matrix - n 6 / weighted by relative Sanderson electronegativities

SpMin6_Bhi

Smallest absolute eigenvalue of Burden modified matrix - n 6 / weighted by relative first ionization potential

SRW9

Self-returning walk count of order 9 ($\ln(1+x)$)

n5Ring

Number of 5-membered rings

AATS5i

Average Broto-Moreau autocorrelation - lag 5 / weighted by first ionization potential

SpMin7_Bhs

Smallest absolute eigenvalue of Burden modified matrix - n 7 / weighted by relative I-state

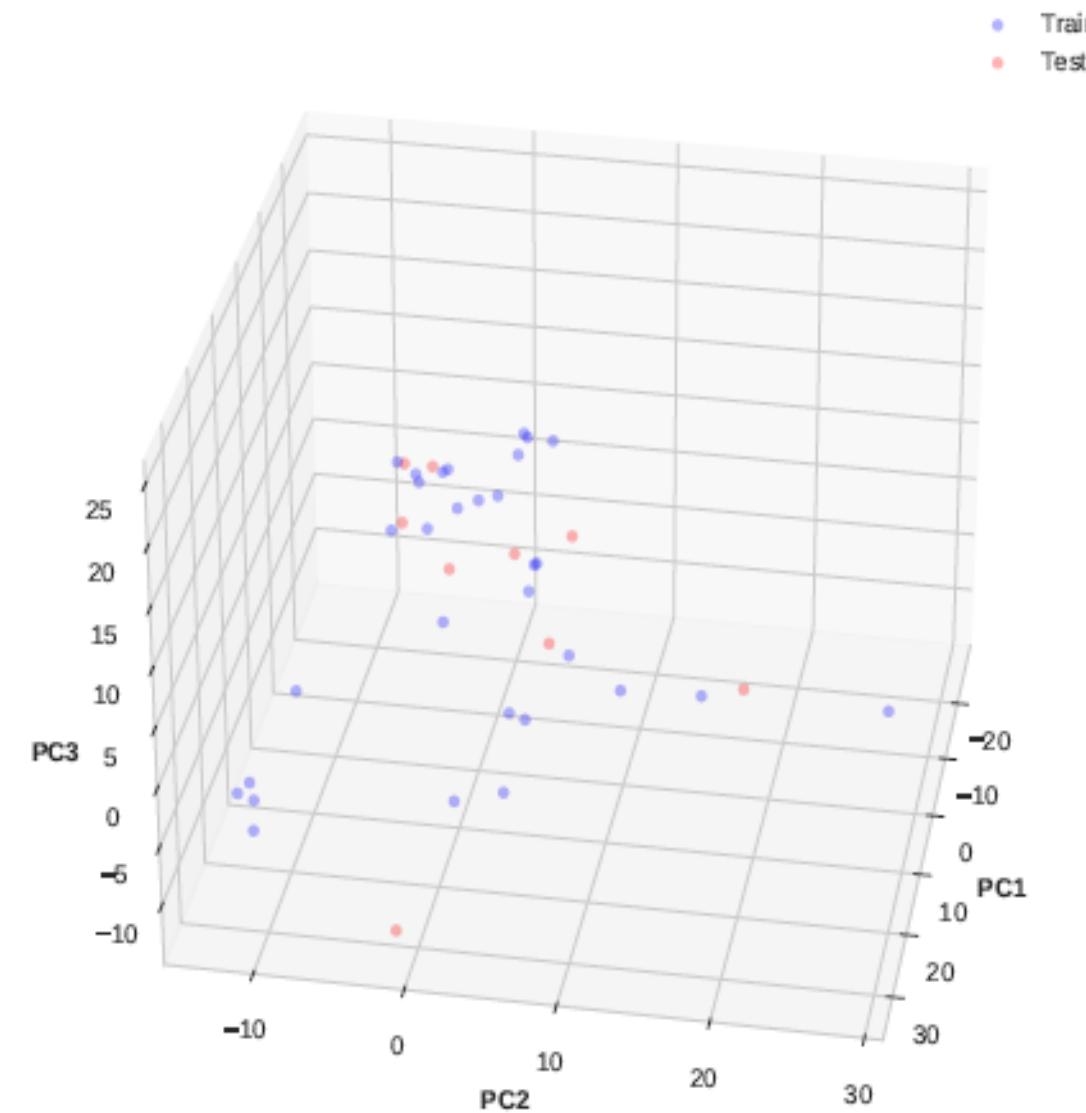
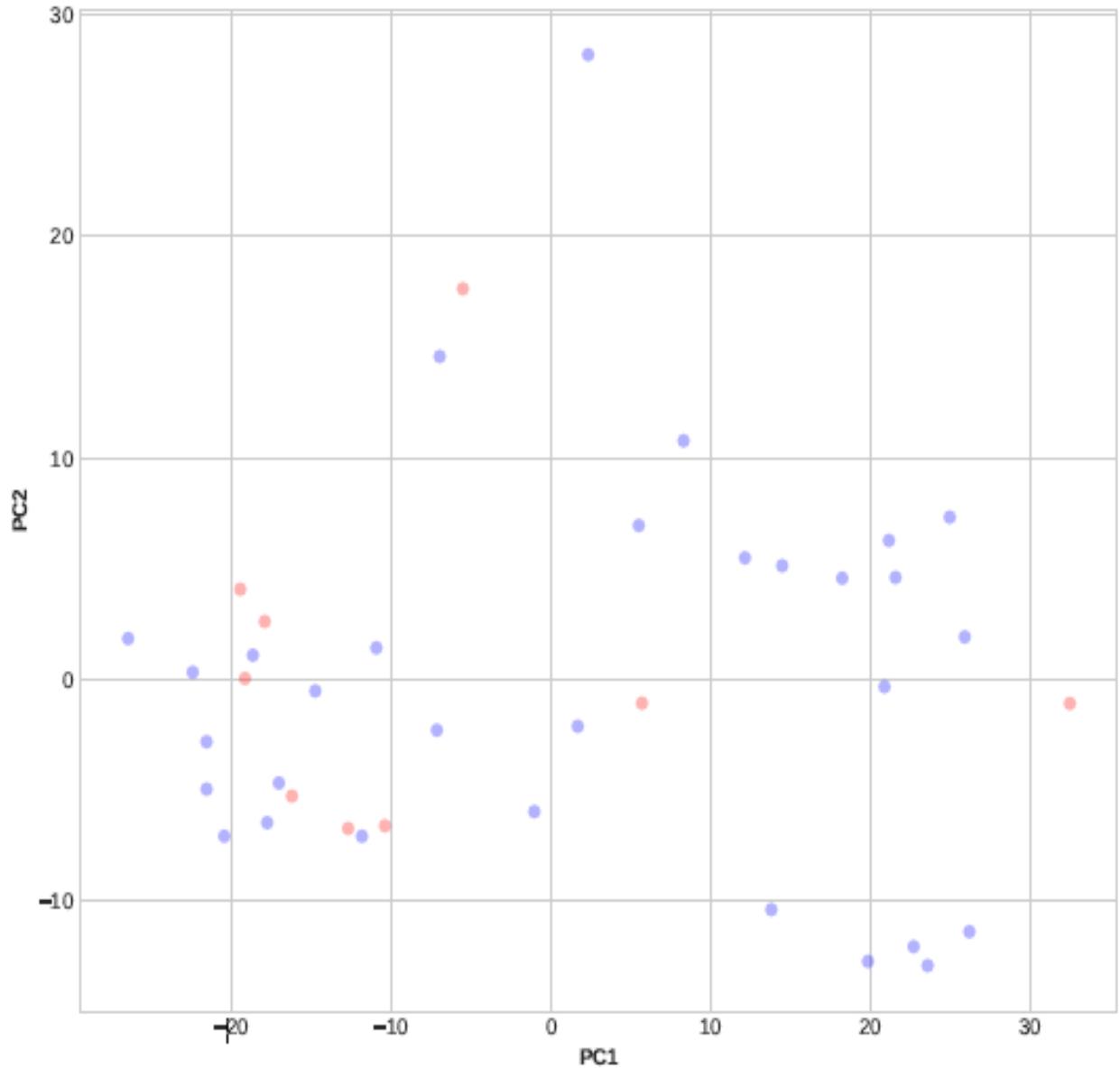
VP-4

Valence path, order 4

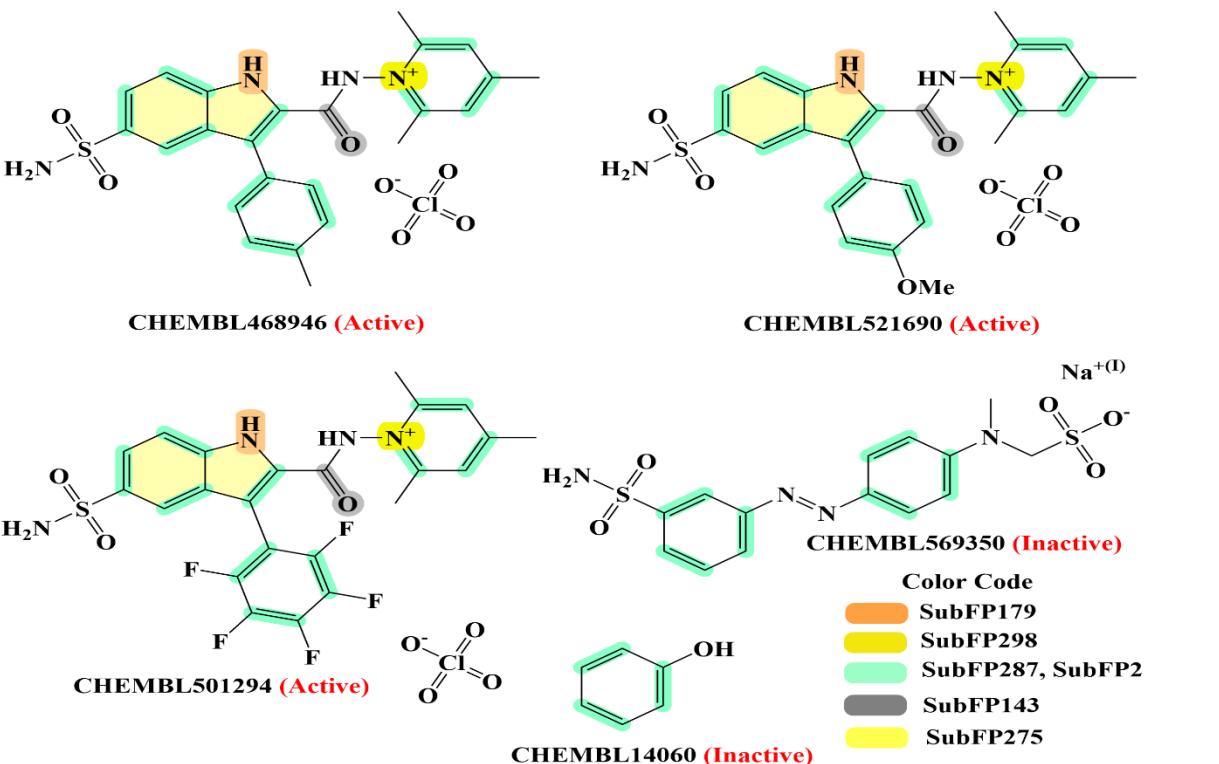
AATSC8i

Average centered Broto-Moreau autocorrelation - lag 8 / weighted by first ionization potential

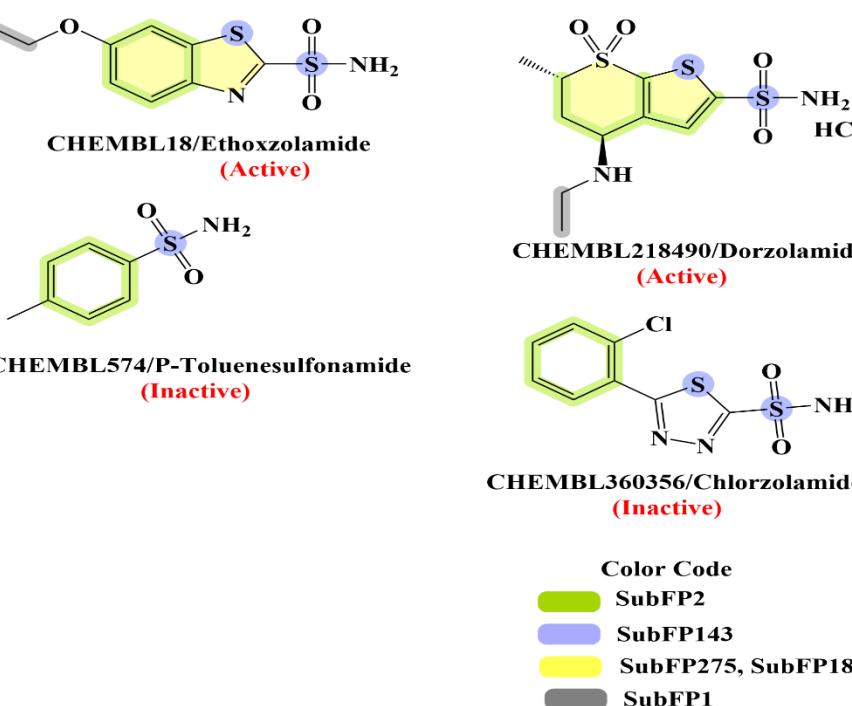
Applicability domain analysis through 2D and 3D PCA plots for *MtbCA2* 1D 2D molecular descriptor prediction model



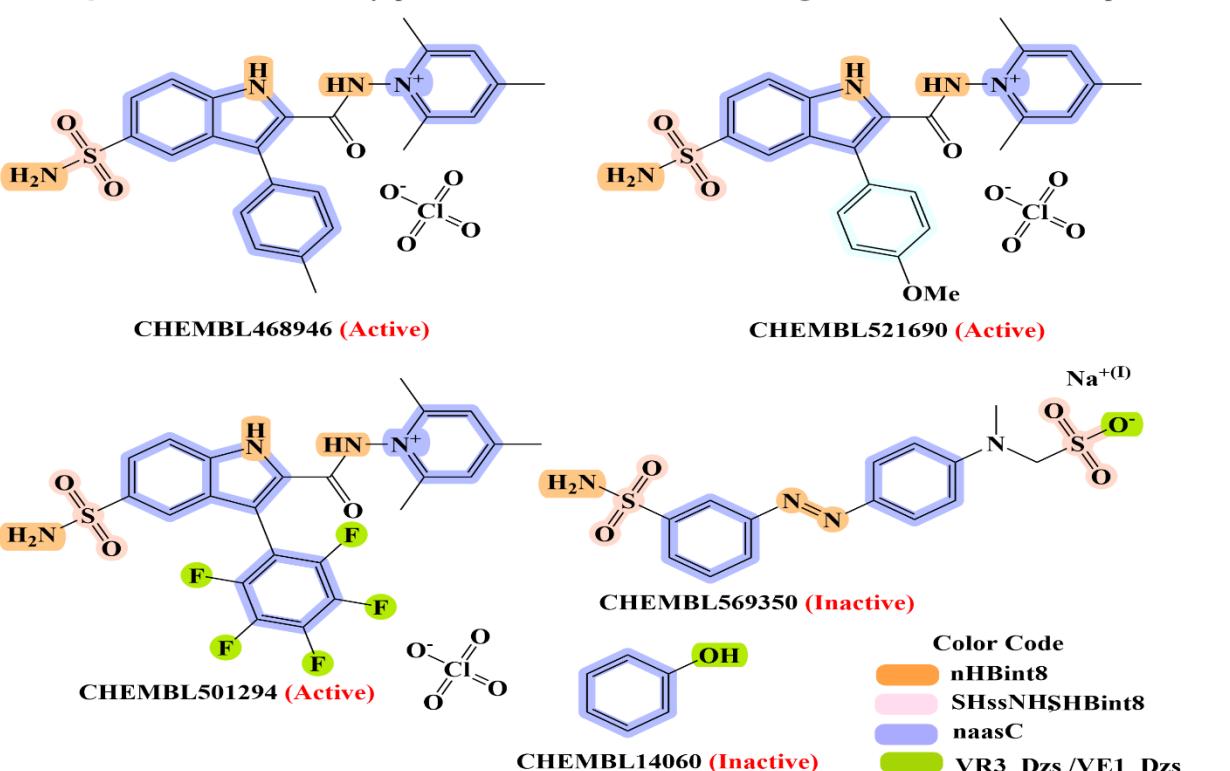
ML-QSAR model for bioactivity prediction of MtbCA1 inhibitors using substructure fingerprints



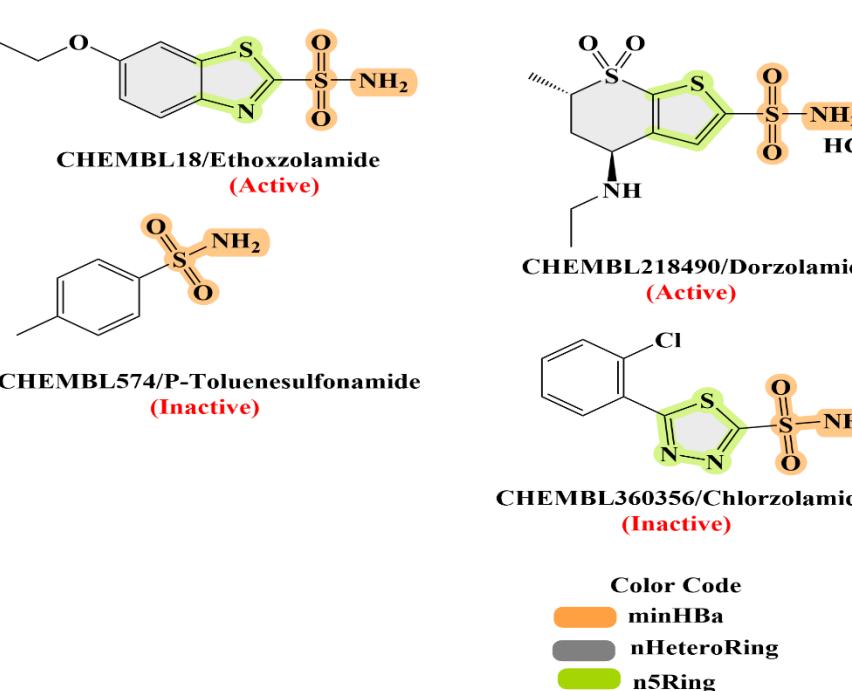
ML-QSAR model for bioactivity prediction of MtbCA2 inhibitors using substructure fingerprints



ML-QSAR model for bioactivity prediction of MtbCA1 inhibitors using 1D and 2D molecular descriptors



ML-QSAR model for bioactivity prediction of MtbCA2 inhibitors using 1D and 2D molecular descriptors



Initial ML-QSAR models deployment using two different molecular signatures with respect to MtbCA1 inhibitors with their respective pKi values)
(S: 307 ; D: 1444)

ML-QSAR models generation and statistical analysis
(S:18; D:763)

Machine Learning models generation with significant molecular signatures

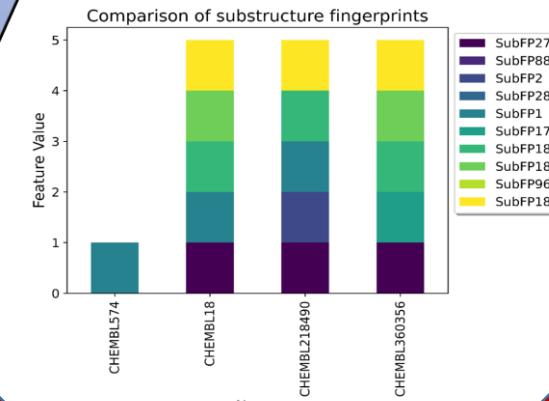
Web app development

Initial ML-QSAR models deployment using two different molecular signatures with respect to MtbCA2 inhibitors with their respective pKi values)
(S: 307 ; D: 1444)

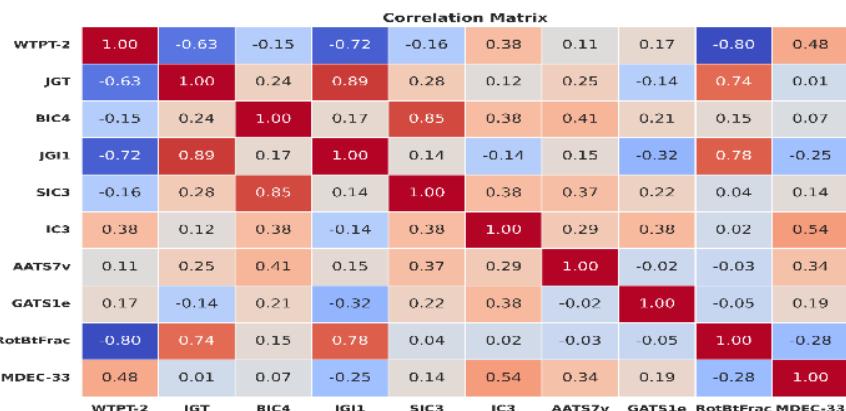
ML-QSAR models generation and statistical analysis
(S:13; D:736)

Machine Learning models generation with significant molecular signatures

Web app development



ML-QSAAR (Dual inhibitors)



MtbCA-Selec-Pred



<https://mtbca-selec-pred.streamlit.app/>

Machine Learning-assisted web application for selectivity bioactivity prediction for MtbCAIs, against MtbCA1 and MtbCA2

S: Substructure Fingerprints D: 1D & 2D Descriptors

X

Choose a prediction model

MtbCA1 prediction model using ...

1. Upload your CSV data

Upload your input file

Drag and drop file here

Limit 200MB per file • TXT

Browse files

[Example input file](#)[Predict](#)

MtbCA-Selec-Pred app

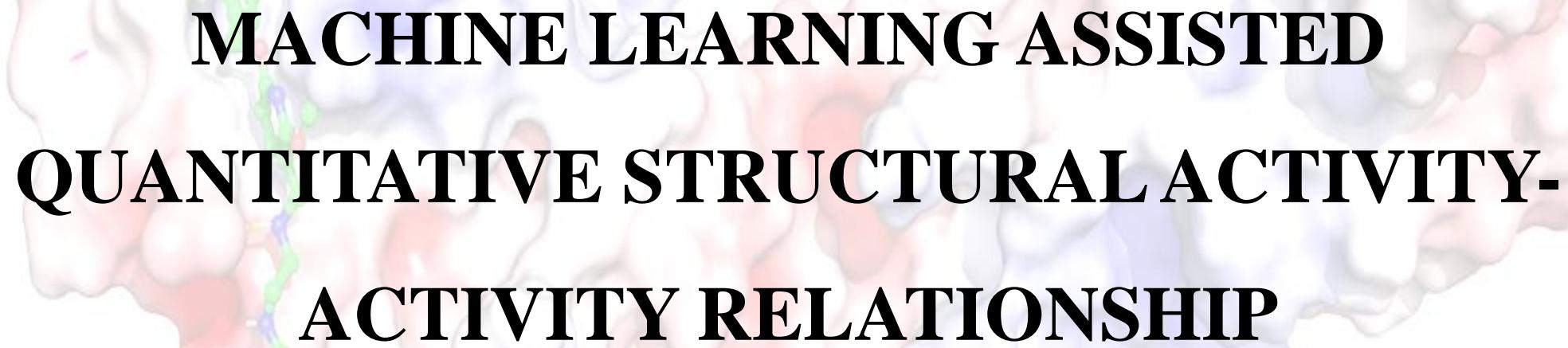
MtbCA-Selec-Pred allows users to predict bioactivity and selectivity of a query molecule separately against the Mycobacterium tuberculosis carbonic anhydrase target protein isoforms, MtbCA1 and MtbCA2

[Main](#) [About](#) [What is Mycobacterium tuberculosis carbonic anhydrase \(MtbCA\)?](#) [Dataset](#) [Model performance](#) [PyTorch](#)

Application Description

This module of [MtbCA-Selec-Pred](#) has been built to predict bioactivity and identify potent inhibitors against Mycobacterium tuberculosis carbonic anhydrases, MtbCA1 and MtbCA2 using robust machine learning algorithms.

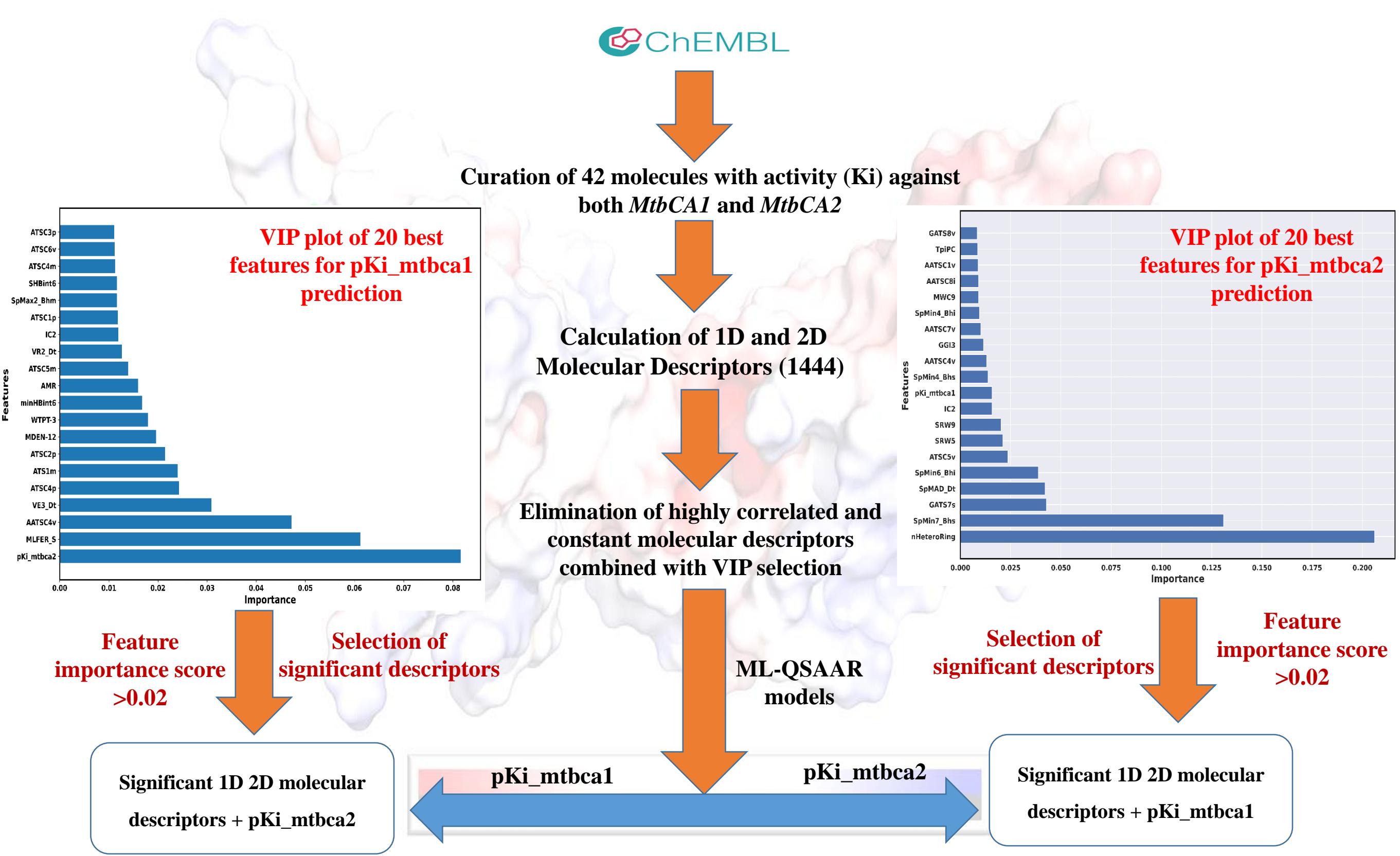
**Predict bioactivity of molecules
against MtbCA1 using substructure
fingerprints**



MACHINE LEARNING ASSISTED QUANTITATIVE STRUCTURAL ACTIVITY- ACTIVITY RELATIONSHIP

-5.000

5.000



ML-QSAR model for bioactivity prediction of MtbCA1 inhibitors using 1D and 2D molecular descriptors



Curation of 42 (MtbCA1 + MtbCA2) inhibitors with Ki value

Total inhibitors:
42

1D and 2D Molecular Descriptors: 1444

Prediction Parameter: Bioactivity Class ('pKi')

Dataset Splitting into training and test set

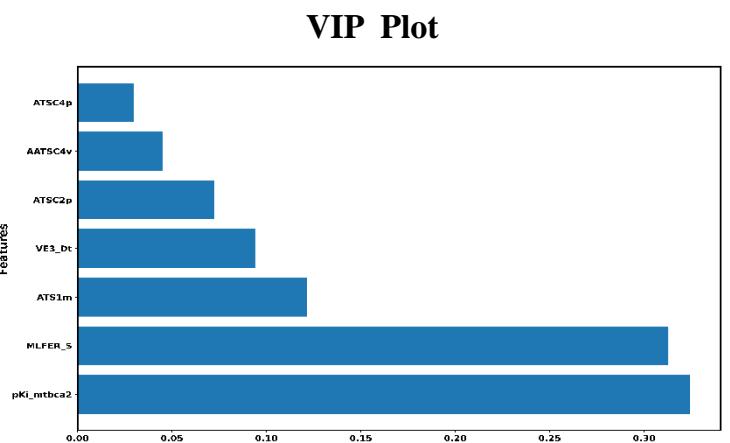
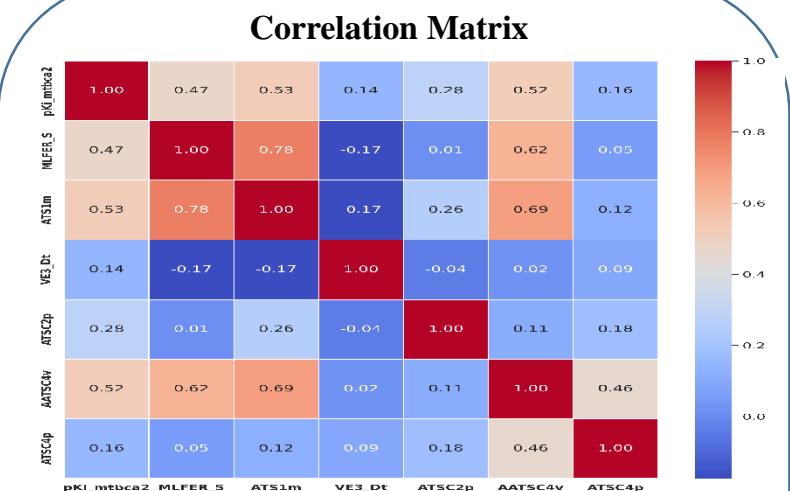
Elimination of highly correlated and constant molecular descriptors

Final number of descriptors(7), number of molecules in training(33) and test(9)

MLP Regressor

Initial Model Performance

Model Metrics
 R2 (Train: 0.9187, Test: 0.4816)
 RMSE (Train: 0.2439, Test: 0.6908)
 MAE (Train: 0.1691, Test: 0.5028)



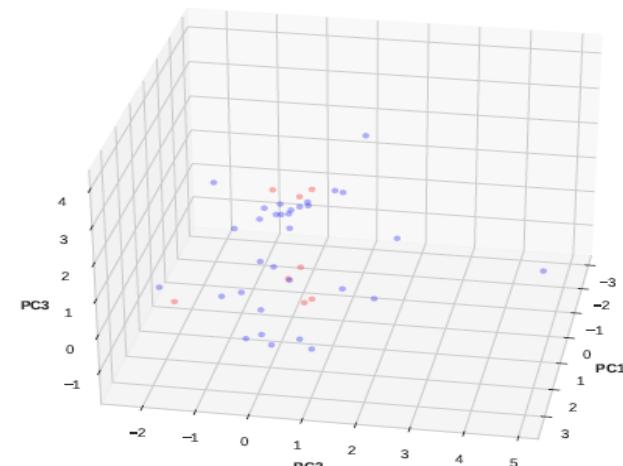
Validation Metrics

Model Metrics
 R2 (Train: 0.9789, Test: 0.794)
 RMSE (Train: 0.122, Test: 0.4878)
 MAE (Train: 0.0904, Test: 0.3647)

Final Model Performance

The final machine-learning model

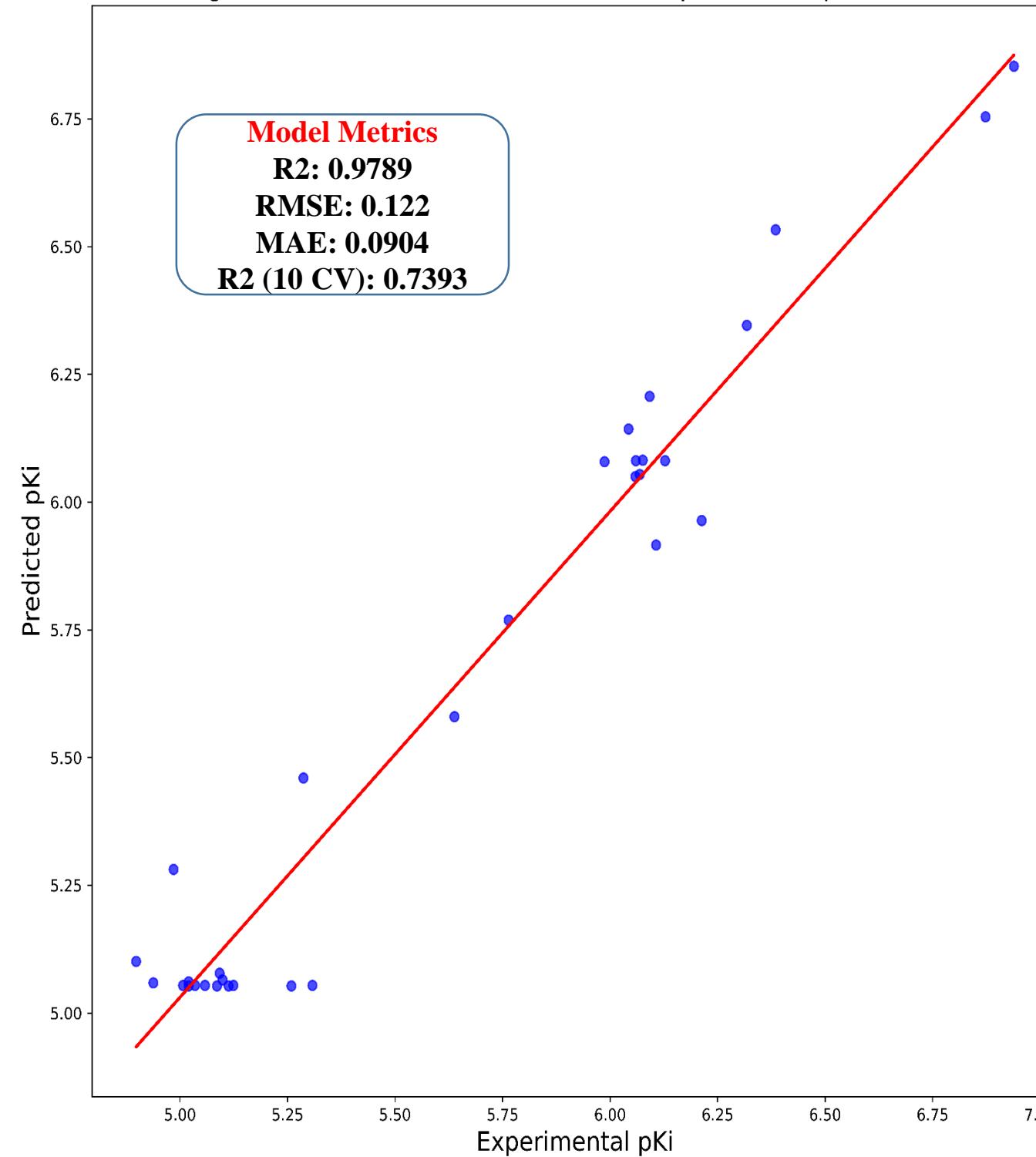
number of molecules in training(32) and test(8)



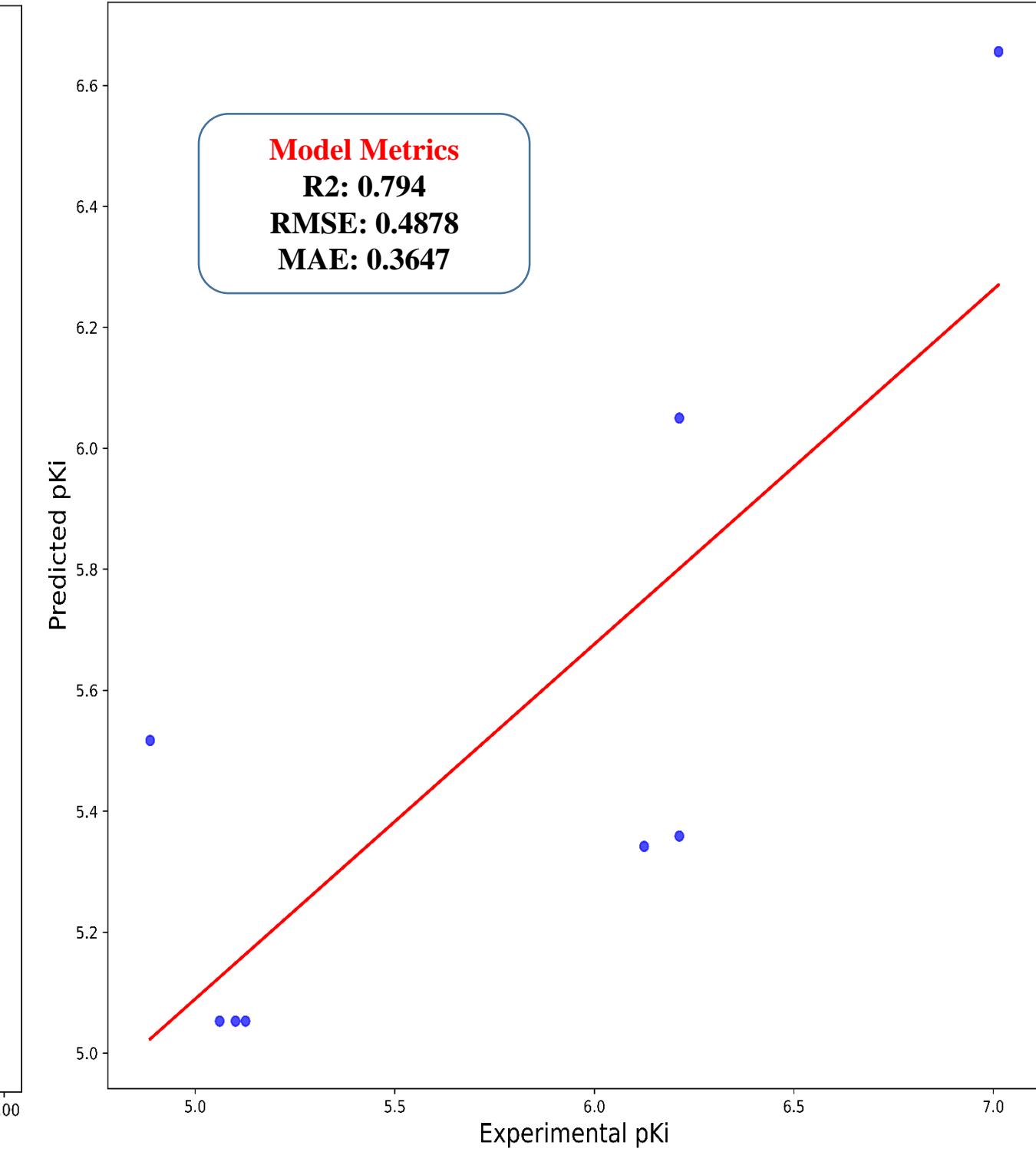
FP and FN Removal

Mention molecules removed (2)

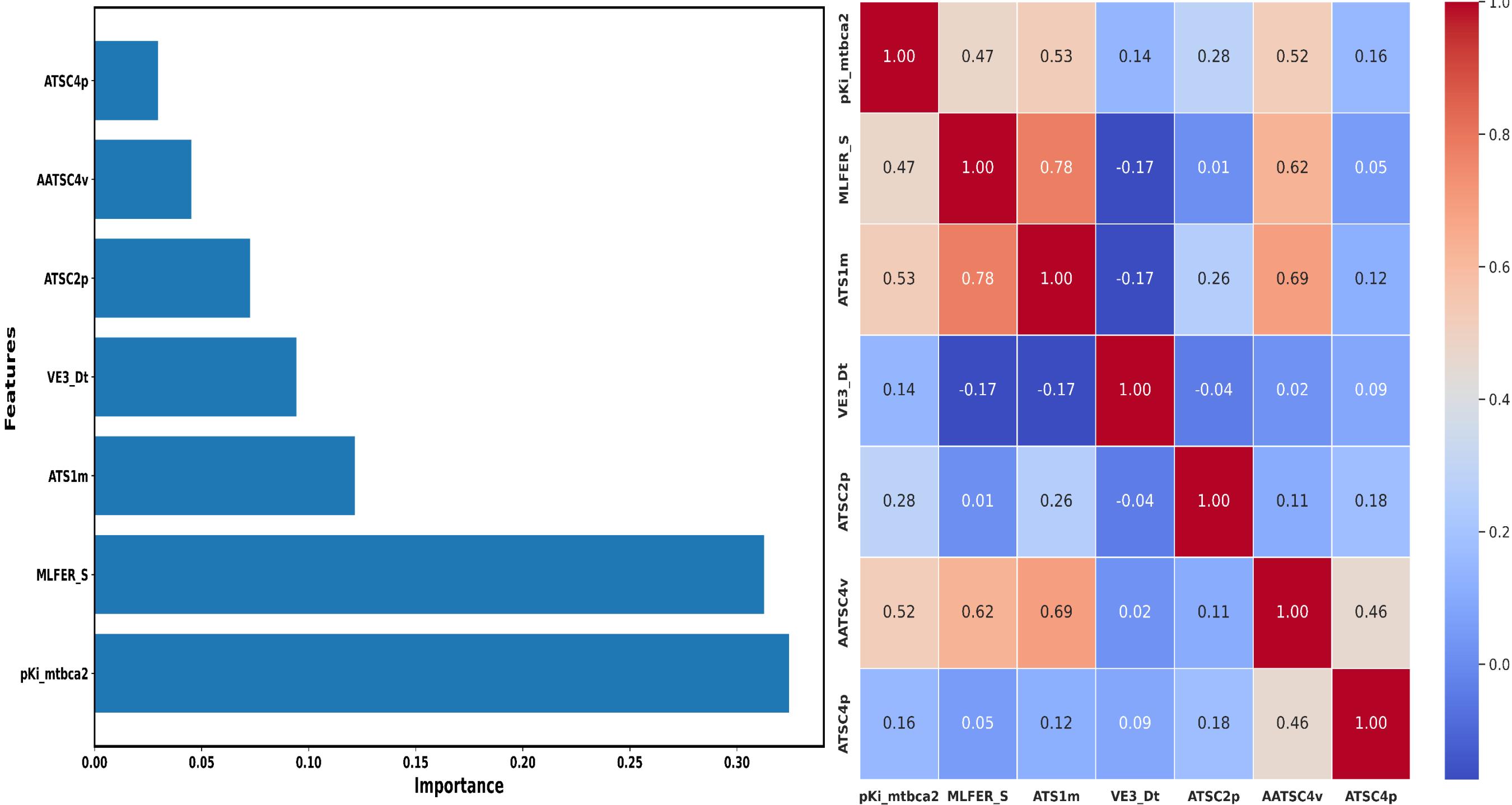
Training set for MtbCA1 1D 2D molecular descriptor QSAAR prediction model



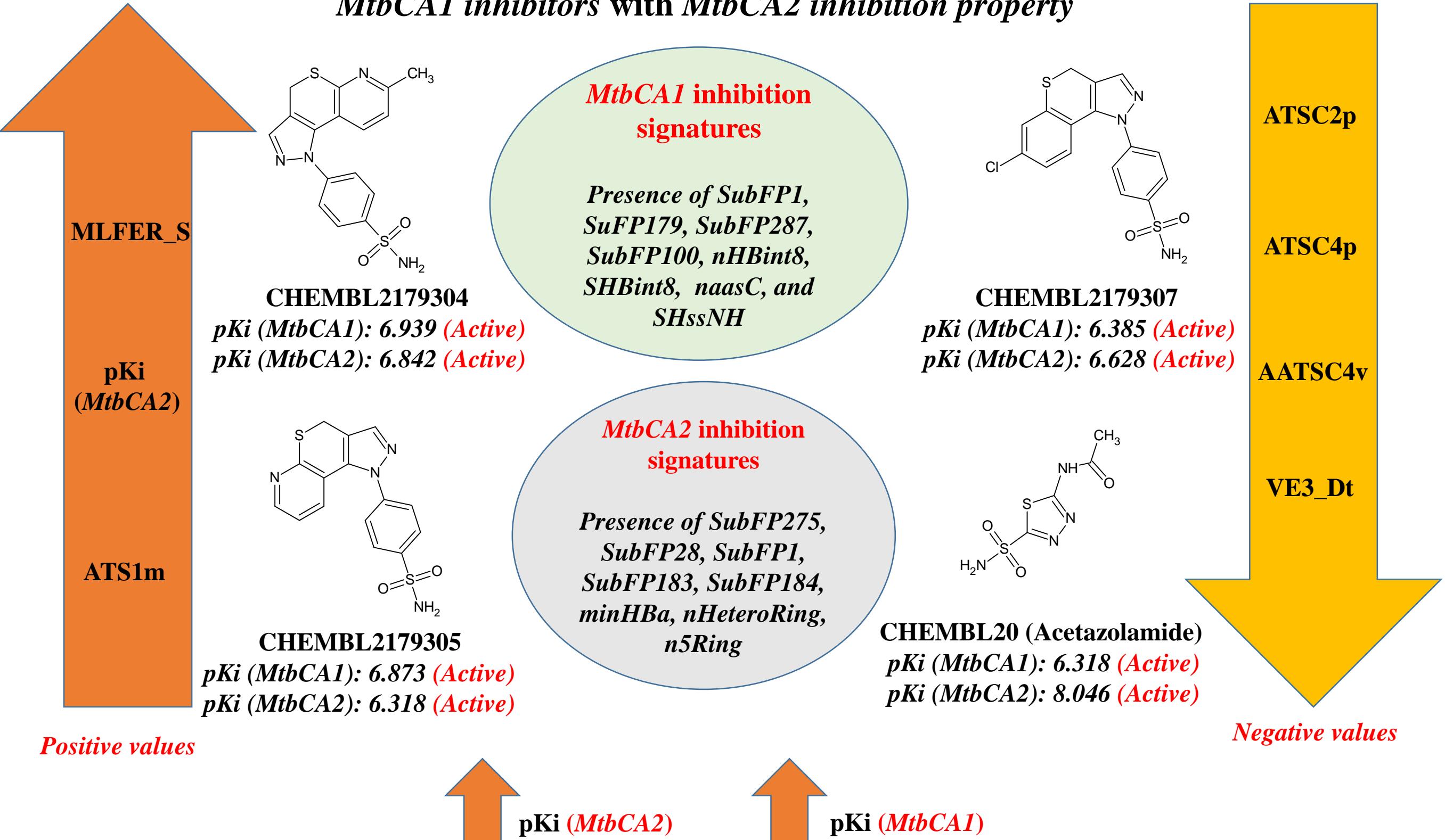
Test set for MtbCA1 1D 2D molecular descriptor QSAAR prediction model



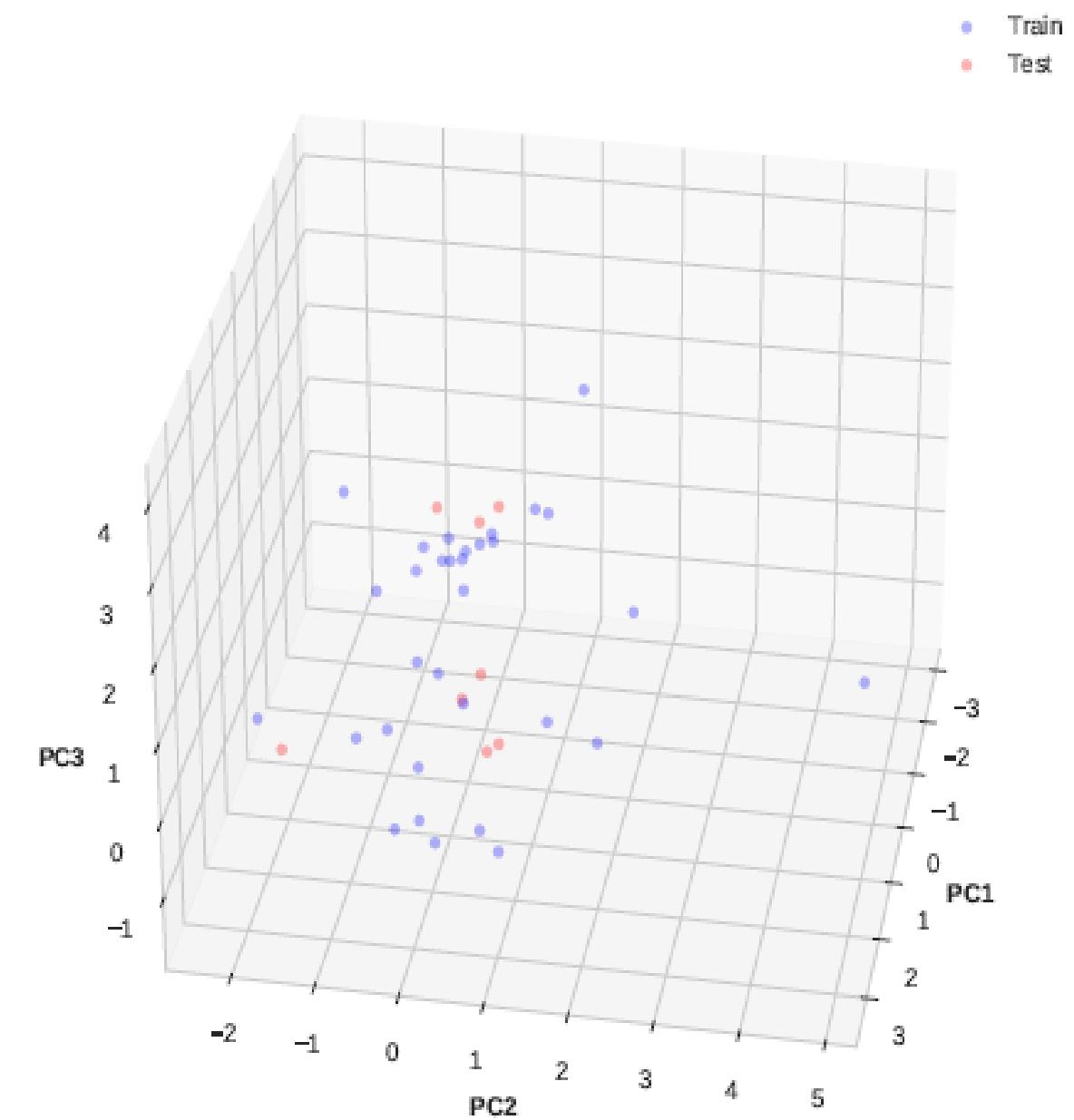
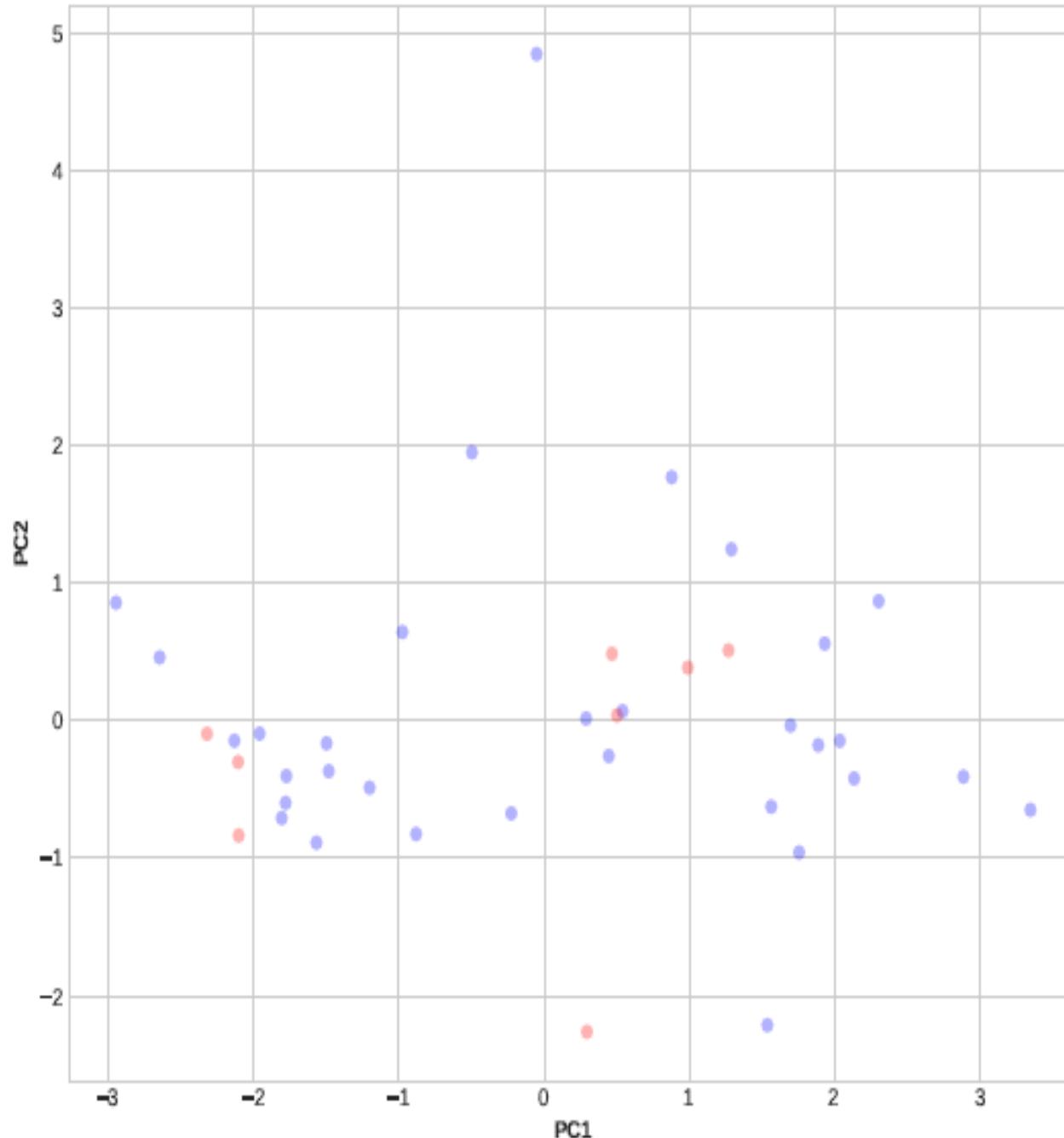
VIP plot and correlation matrix analysis for *MtbCA1* 1D 2D molecular descriptor QSAAR prediction model



MtbCA1 inhibitors with *MtbCA2* inhibition property



Applicability domain analysis through 2D and 3D PCA plots for *MtbCA1* 1D 2D molecular descriptor QSAAR prediction model



ML-QSAAR model for bioactivity prediction of *MtbCA2* inhibitors using 1D and 2D molecular descriptors



Curation of 42 (*MtbCA1 + MtbCA2*) inhibitors with K_i value

Total inhibitors:
42

1D and 2D Molecular Descriptors: 1444

Prediction Parameter: Bioactivity Class ('pKi')

Dataset Splitting into training and test set

Elimination of highly correlated and constant molecular descriptors

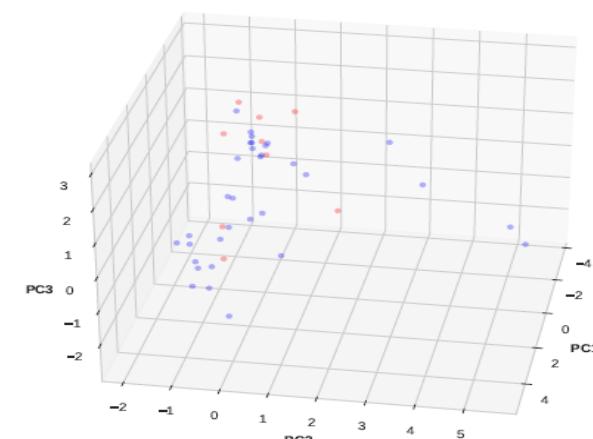
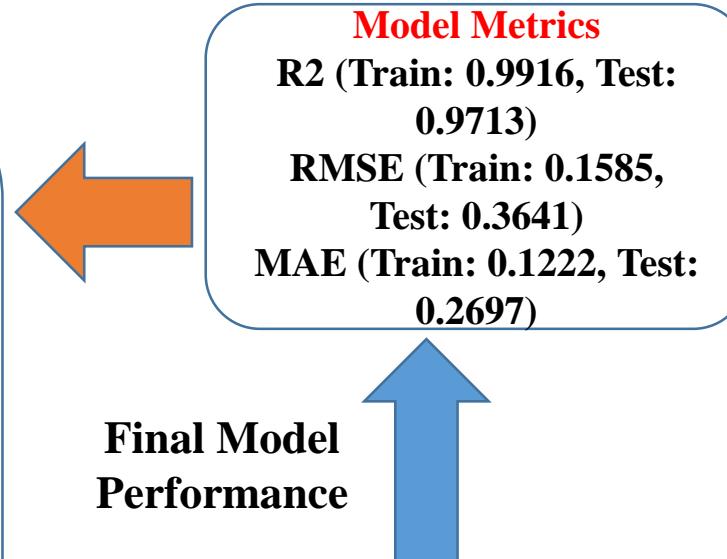
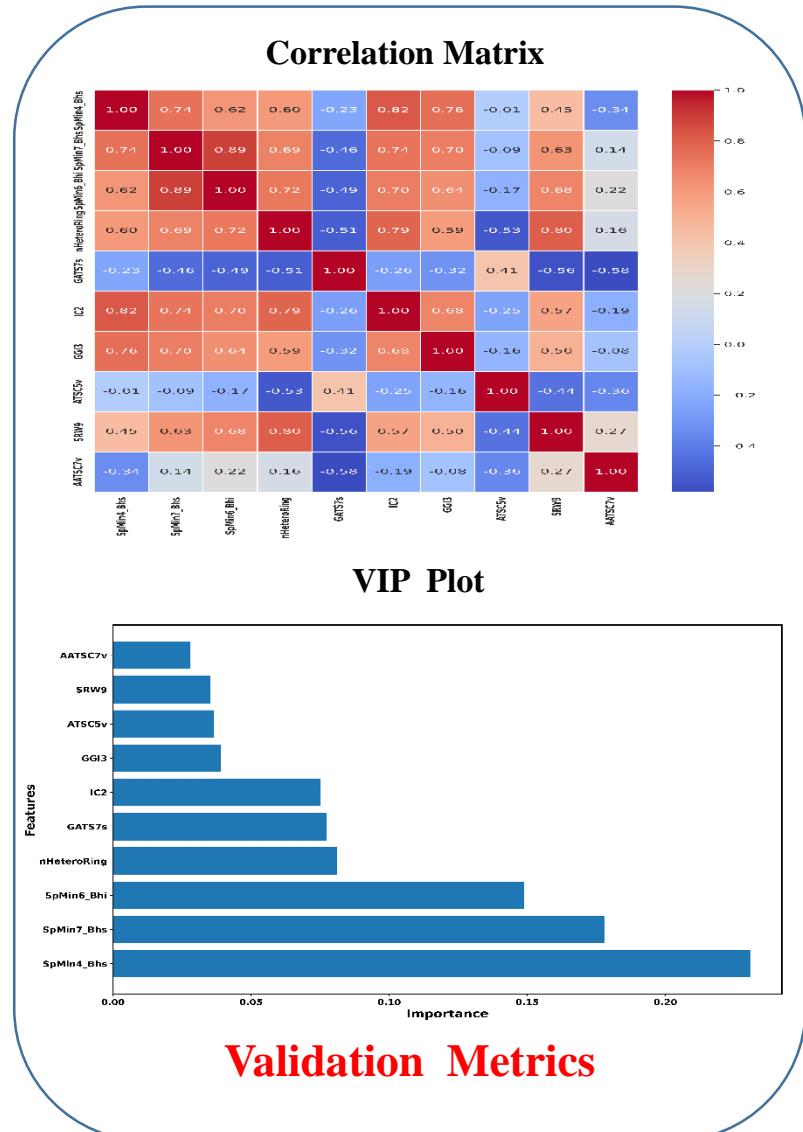
Final number of descriptors(7), number of molecules in training(33) and test(9)

MLP Regressor
Initial Model Performance

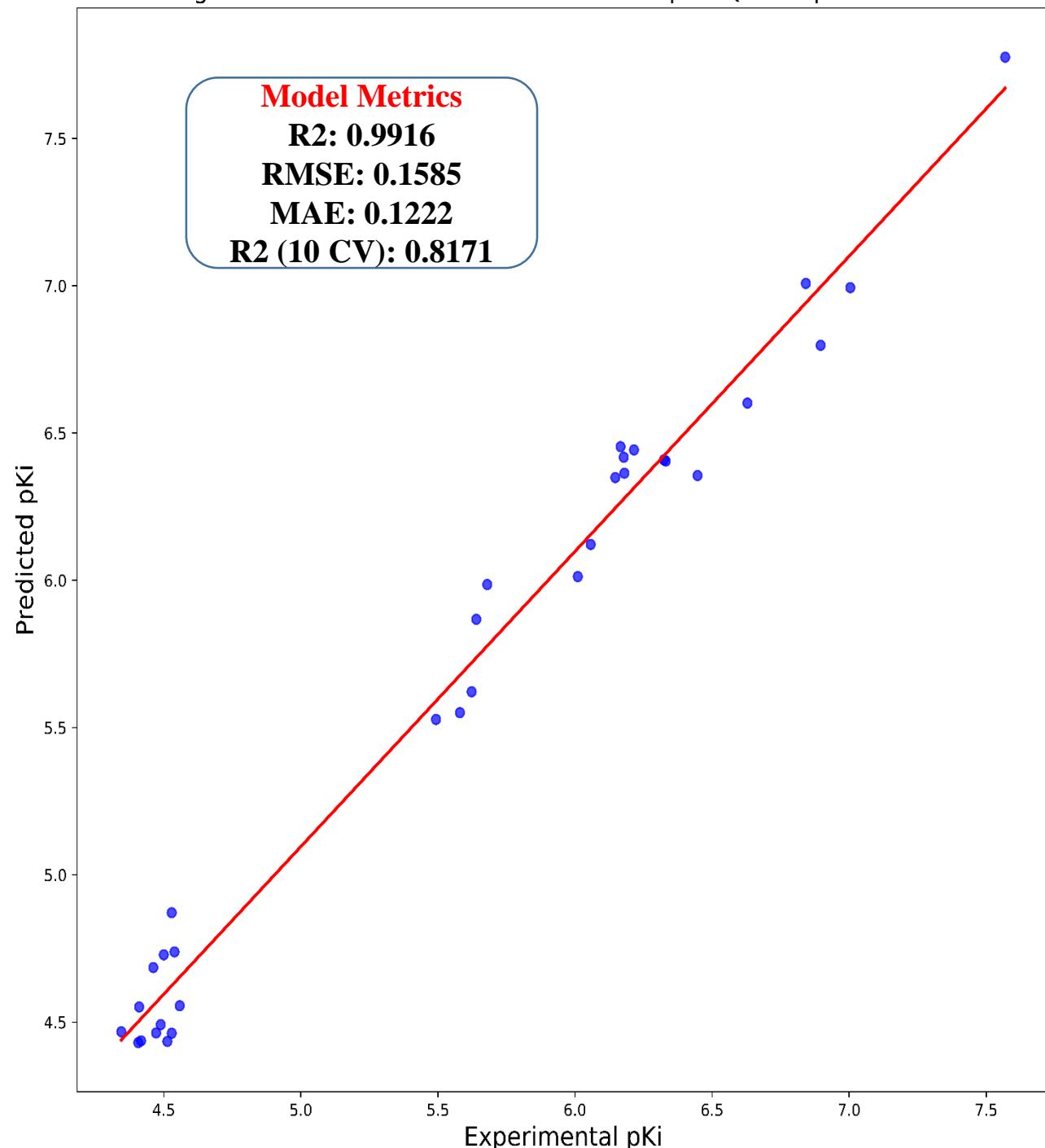
Model Metrics
R2 (Train: 0.9916, Test: 0.8612)
RMSE (Train: 0.1585, Test: 0.7041)
MAE (Train: 0.1222, Test: 0.5742)

FP and FN Removal
Mention molecules removed (2)

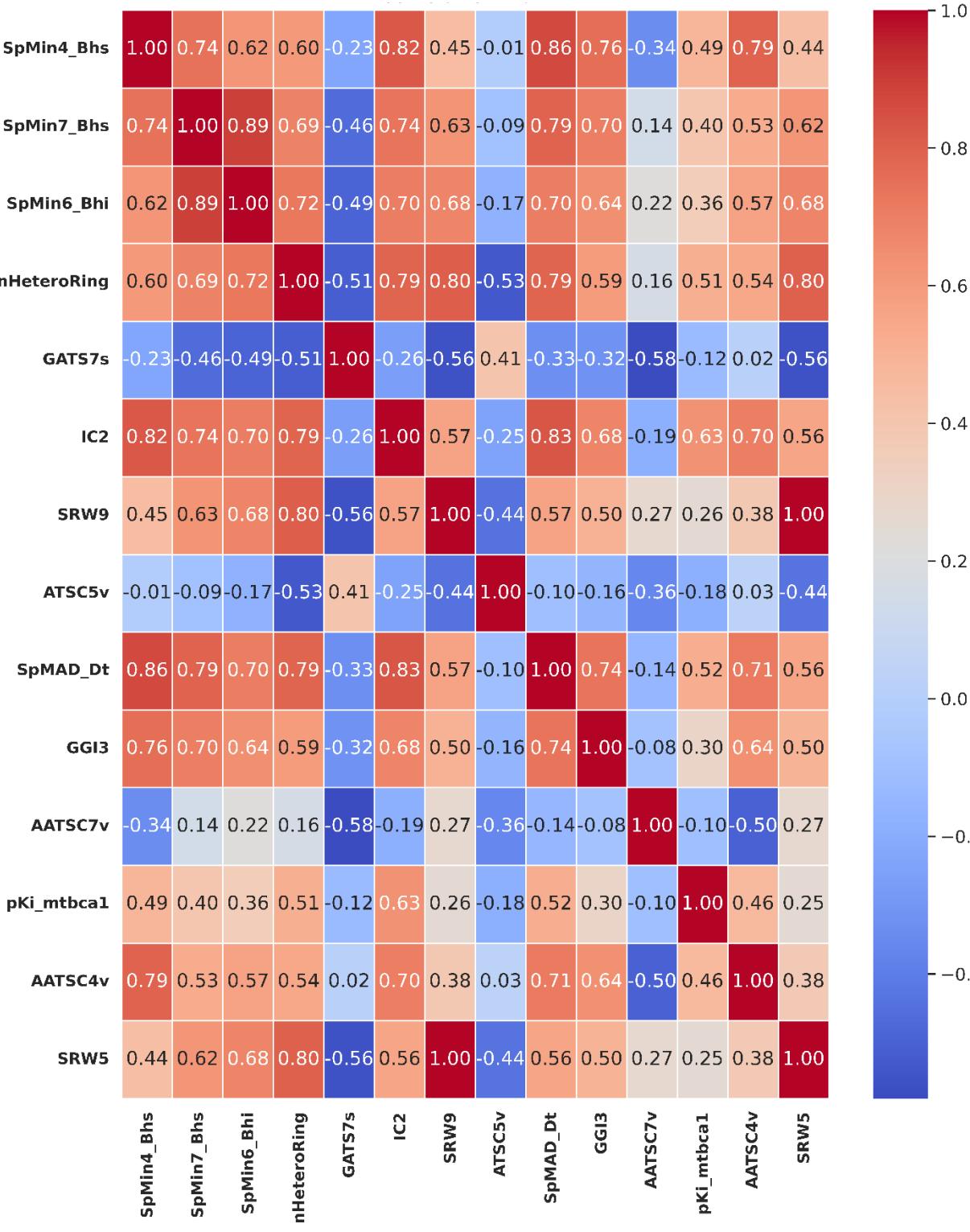
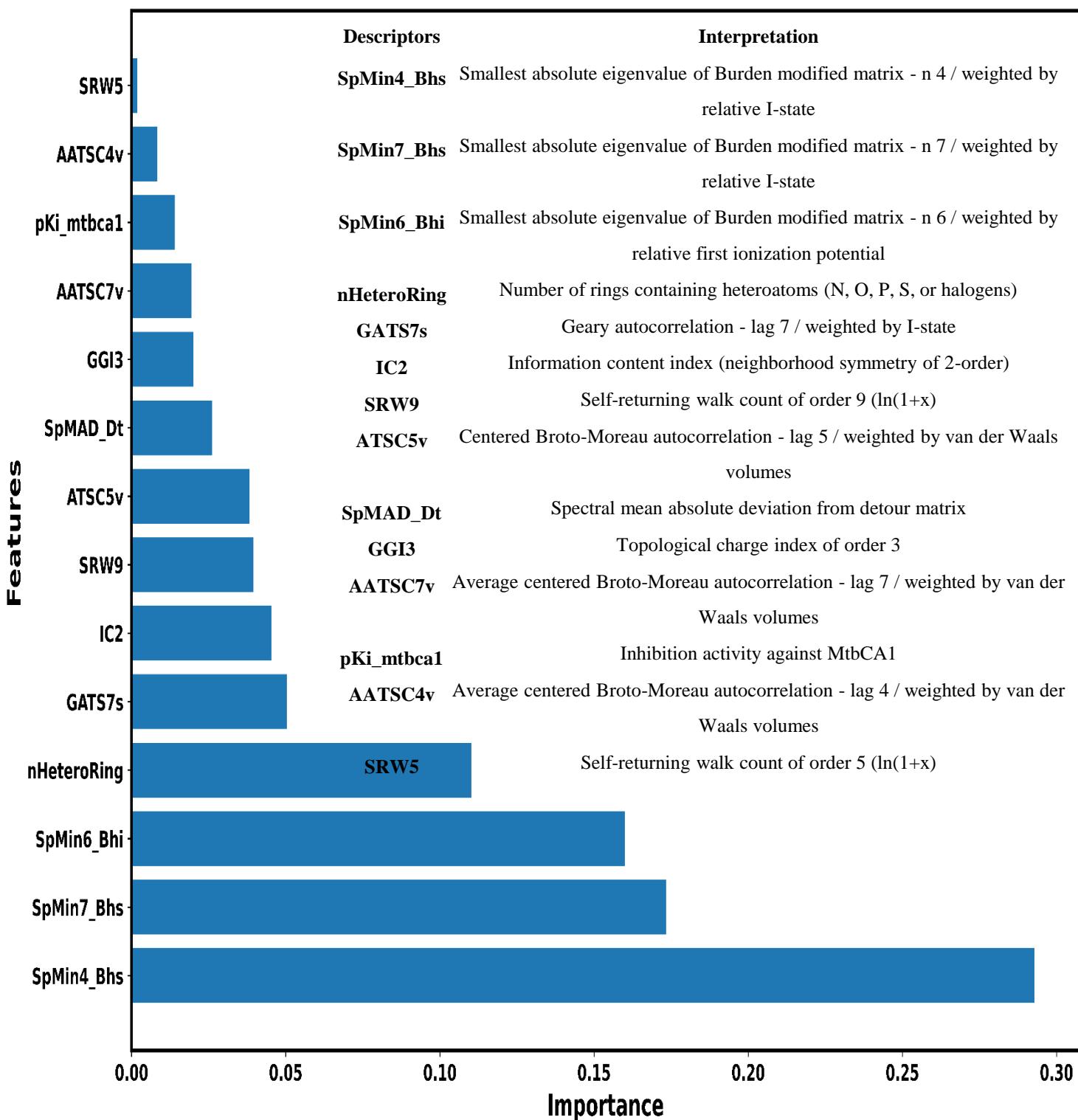
Applicability Domain Analysis



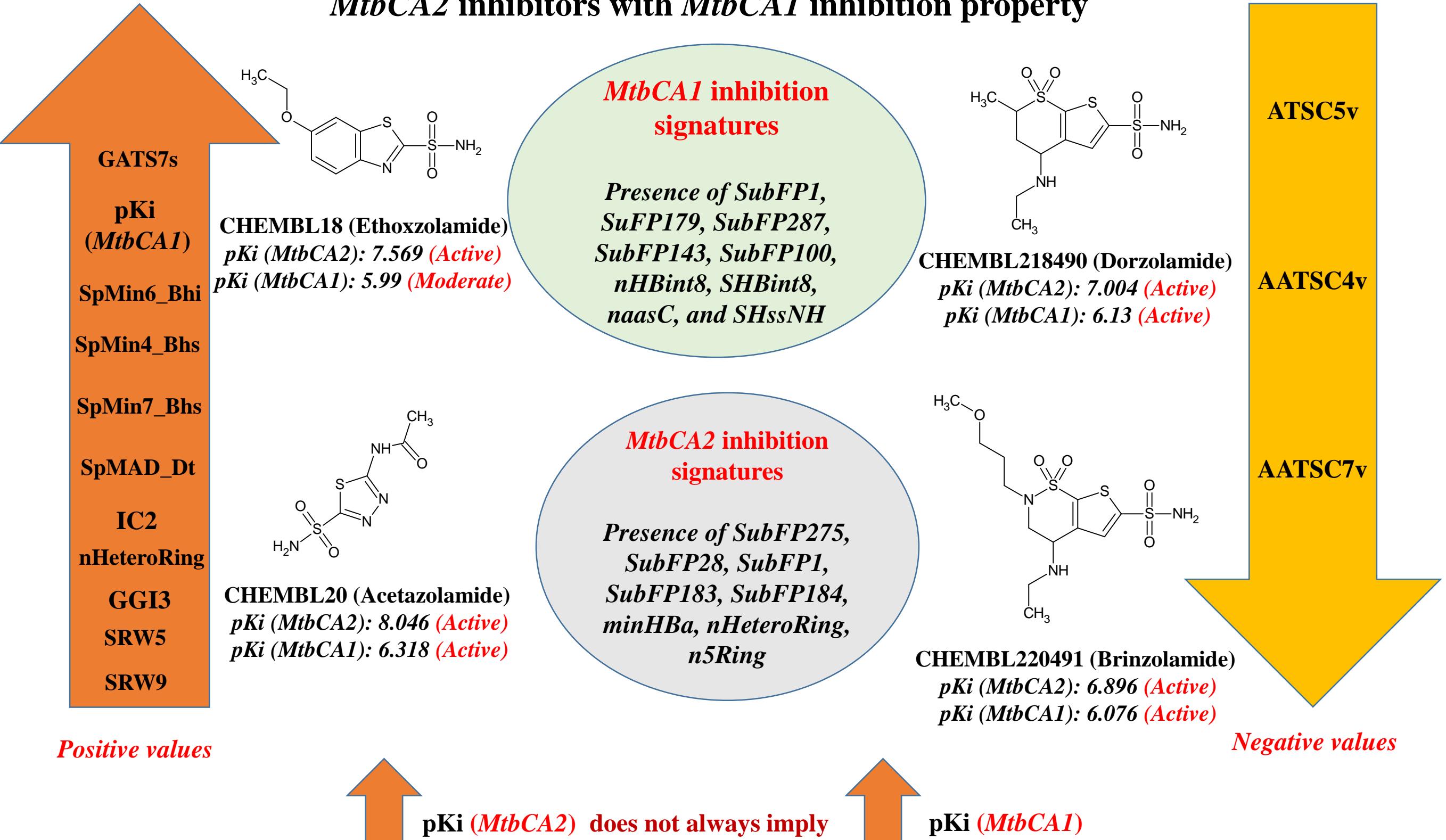
Training set for MtbCA2 1D 2D molecular descriptor QSAAR prediction model



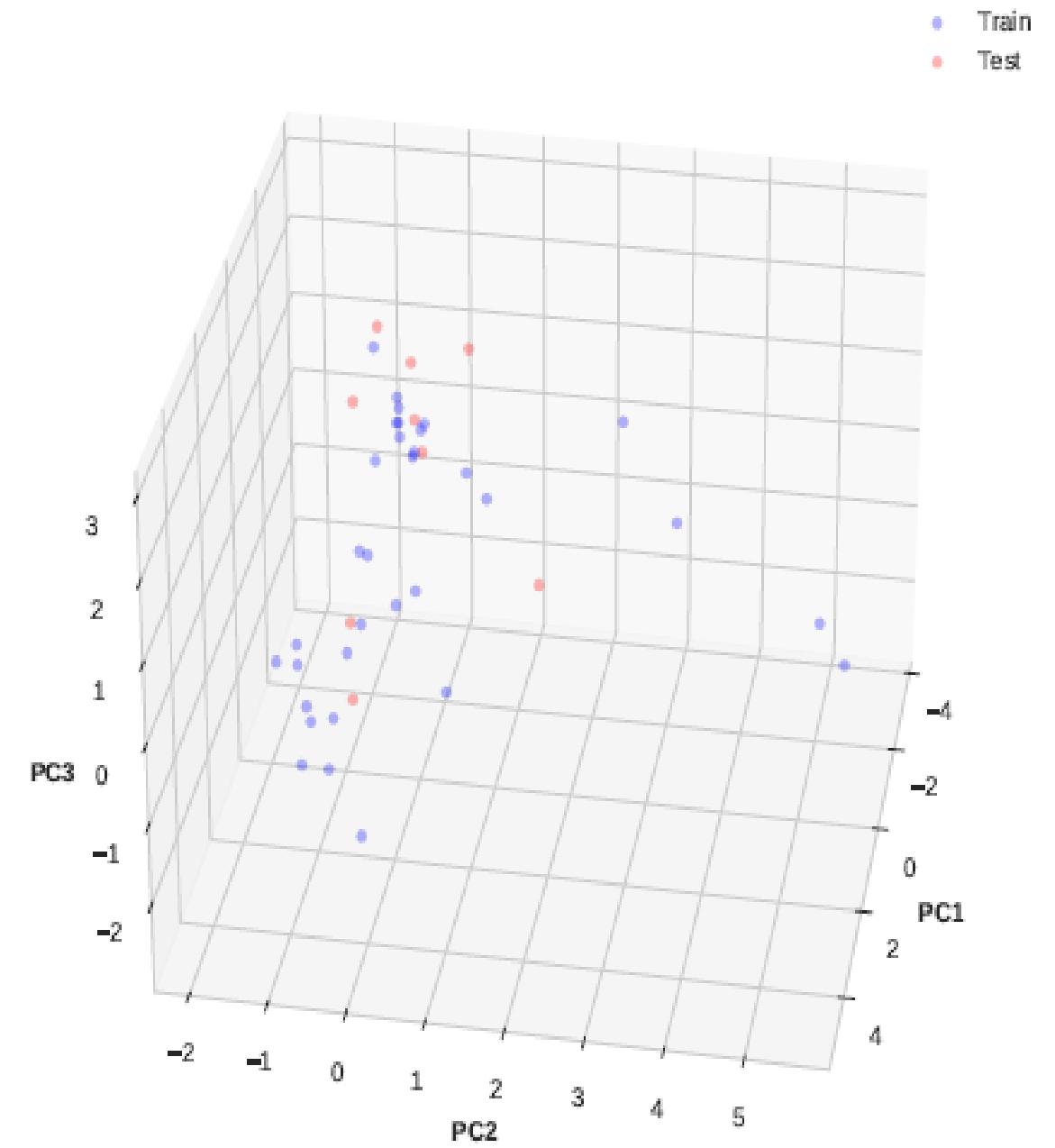
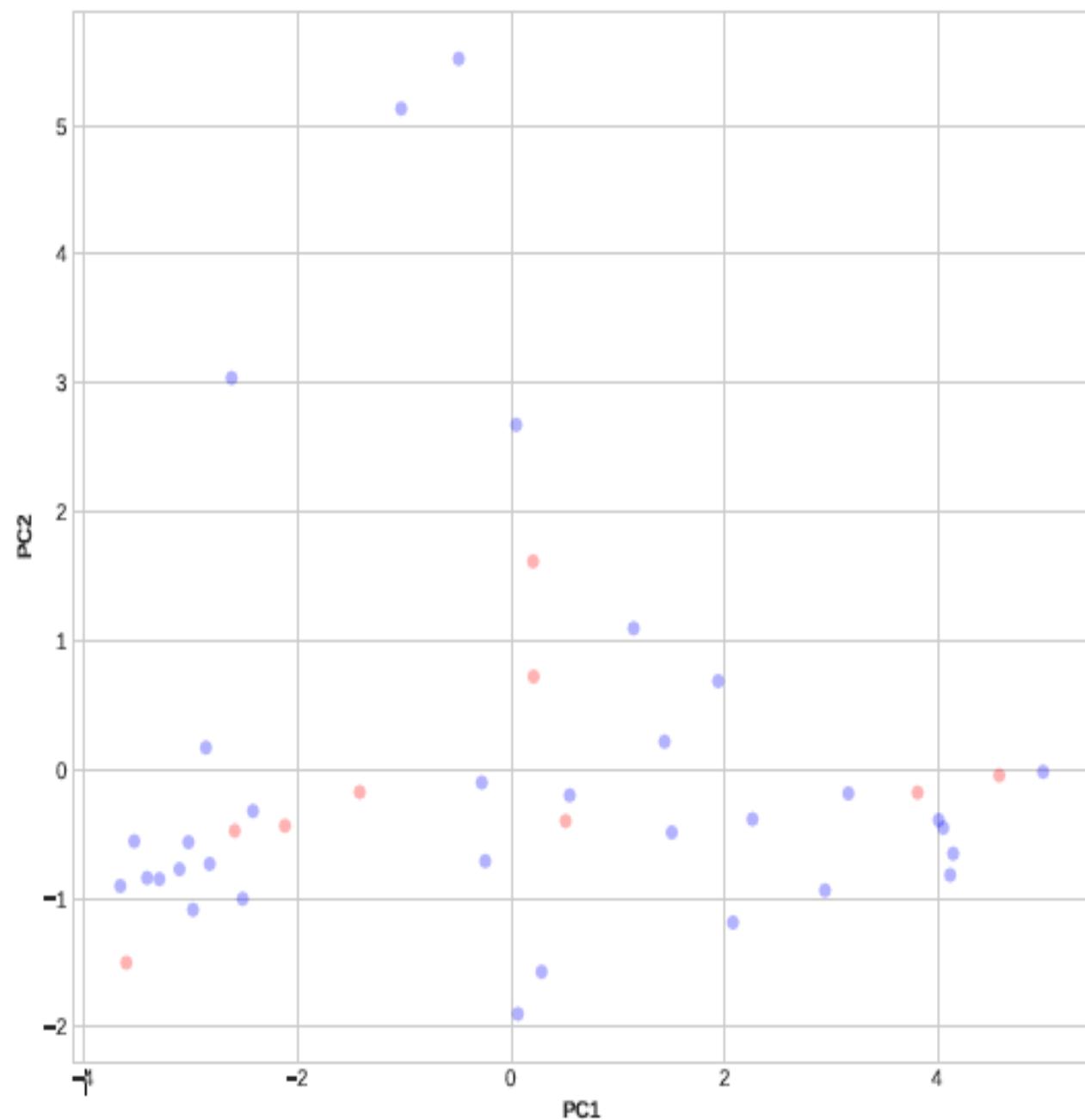
VIP plot and correlation matrix analysis for *MtbCA2* 1D 2D molecular descriptor QSAAR prediction model



MtbCA2 inhibitors with *MtbCA1* inhibition property



Applicability domain analysis through 2D and 3D PCA plots for *MtbCA2* 1D 2D molecular descriptor QSAAR prediction model



ML-QSAAR model for bioactivity prediction of MtbCA1 inhibitors using 1D and 2D molecular descriptors



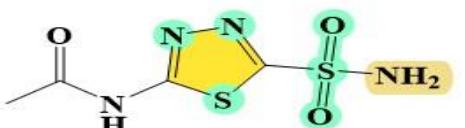
CHEMBL2179304
pKi (MtbCA1): 6.939 (Active)
pKi (MtbCA2): 6.842 (Active)



CHEMBL2179305
pKi (MtbCA1): 6.873 (Active)
pKi (MtbCA2): 6.318 (Active)



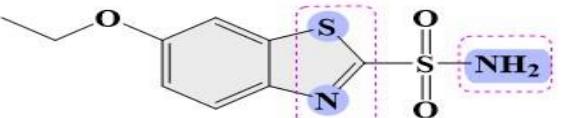
CHEMBL2179307
pKi (MtbCA1): 6.385 (Active)
pKi (MtbCA2): 6.628 (Active)



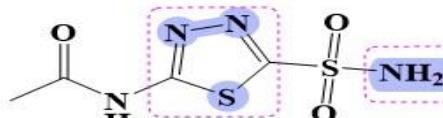
CHEMBL20 (Acetazolamide)
pKi (MtbCA1): 6.318 (Active)
pKi (MtbCA2): 8.046 (Active)

Color Code
█ MLFER_S
█ VE3_Dt
█ ATS1m

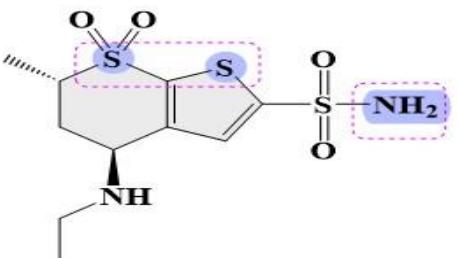
ML-QSAAR model for bioactivity prediction of MtbCA2 inhibitors using 1D and 2D molecular descriptors



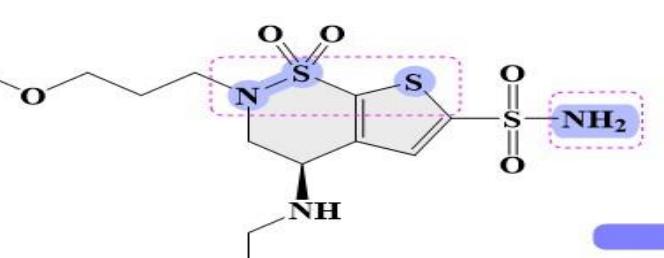
CHEMBL18 (Ethoxzolamide)
pKi (MtbCA2): 7.569 (Active)
pKi (MtbCA1): 5.99 (Moderate)



CHEMBL20 (Acetazolamide)
pKi (MtbCA1): 6.318 (Active)
pKi (MtbCA2): 8.046 (Active)

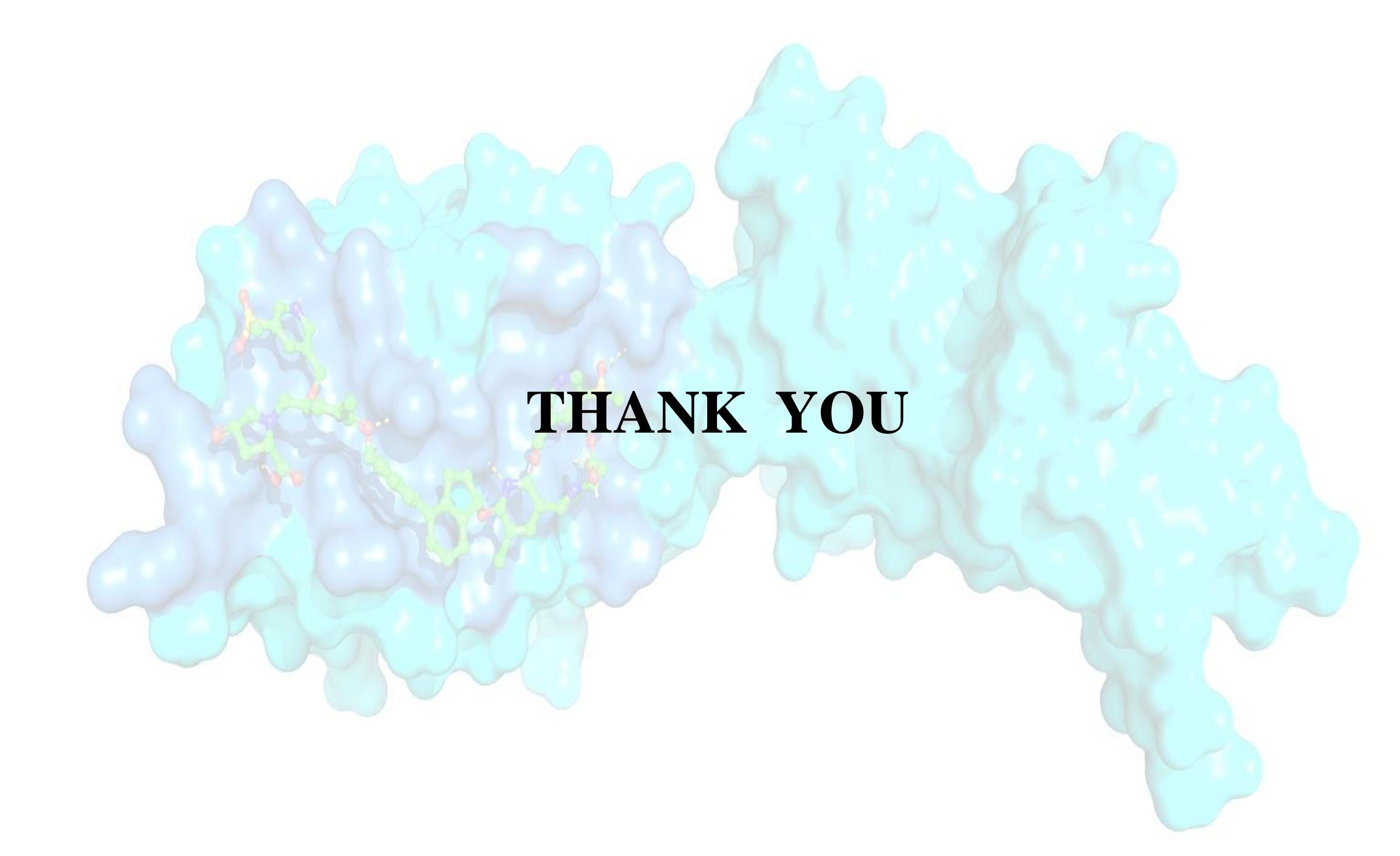


CHEMBL218490 (Dorzolamide)
pKi (MtbCA2): 7.004 (Active)
pKi (MtbCA1): 6.13 (Active)



CHEMBL220491 (Brinzolamide)
pKi (MtbCA2): 6.896 (Active)
pKi (MtbCA1): 6.076 (Active)

Color Code
█ GGI3
█ SpMin4_Bhs,
SpMin7_Bhs, and
SpMin6_Bhs
█ nHeteroRing



THANK YOU