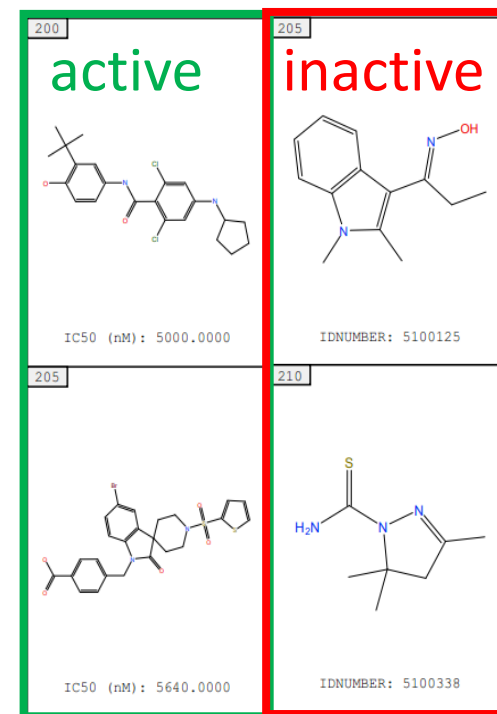
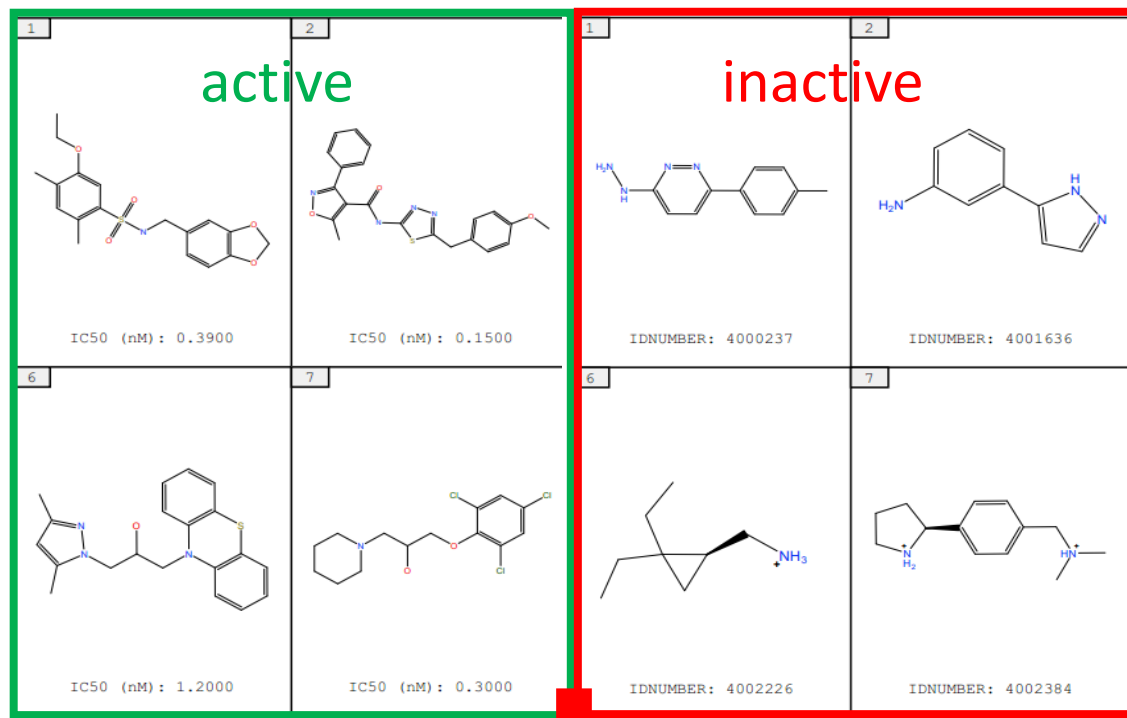


# 基于机器学习的判别模型构建

# 判别模型建模一般流程

- 已知活性数据收集
- 数据集预处理（正样本/负样本、训练集/测试集准备等）
- 分子描述属性计算（传统分子描述符、分子指纹等）
- 模型构建（机器学习算法）
- 外部数据集检测
- 未知化合物类别预测

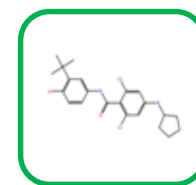
# 以朴素贝叶斯为例对FXR活性剂与非活性剂进行ML判别模型构建



标签值	描述符Xa	描述符Xb	描述符Xc	.....
active	Xa1	Xb1	Xc1	.....
active	Xa2	Xb2	Xc2	.....
inactive	Xa3	Xb3	Xc3	.....
.....	.....	.....	.....	.....

Bayes

描述符Xa	描述符Xb	描述符Xc	.....
Xa1	Xb1	Xc1	.....
Xa2	Xb2	Xc2	.....
Xa3	Xb3	Xc3	.....



active, 80%  
inactive, 20%



# 数据集准备

Small Molecules

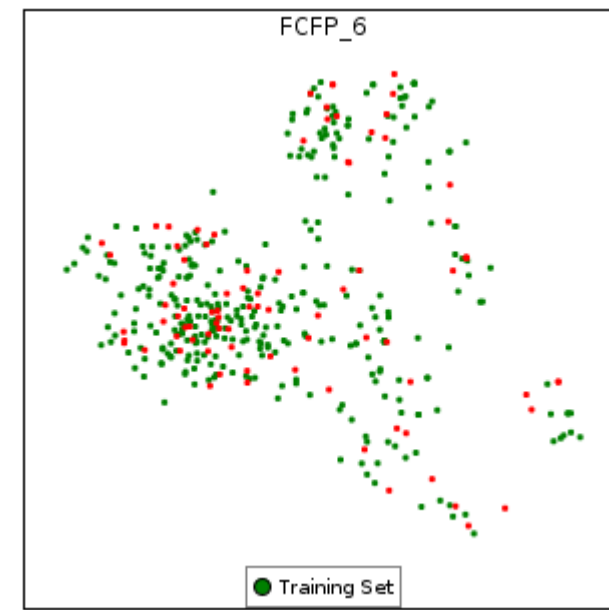


Create QSAR Model



**Prepare Data**

- Prepare Ligands for QSAR...
- Prepare Dependent Property...
- Generate Training and Test Data...**



Generate Training and Test Data

Parameter Name	Parameter Value
Input Ligands	dataset:All
> Split Method	Random
Training Set Percentage	80
> Use Properties	PredefinedSet
> Clustering Options	

☐ Show Parameter Help

Run Options Cancel Help

- TrainingSet
- TestSet

# 模型的参数设置与构建

**Build Models**

Use the following tools to build classification models.

Create Bayesian Model...

Create Recursive Partitioning Model...

Use the following tools to build regression models.

Create Multiple Linear Regression Model...

Create Partial Least Squares Model...

Create Genetic Function Approximation Model...

Use the following tool to build molecular field analysis models.

Create 3D QSAR Model...



	Index	Name	Visible	Tagged	Visibility Locked	Active
1	1	Molecule	<input checked="" type="checkbox"/> Yes	<input type="checkbox"/> No	<input type="checkbox"/> No	0
2	2	Binding...	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1
3	3	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1
4	4	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	0
5	5	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	0
6	6	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	0
7	7	Binding...	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1
8	8	Binding...	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1
9	9	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1
10	10	Binding...	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1



Create Bayesian Model

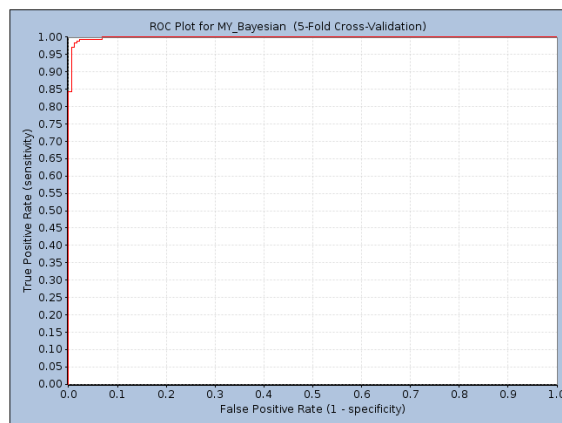
Parameter Name	Parameter Value
Input Ligands	TrainingSet_dataset:All
Input Test Ligands	TestSet_dataset:All
Property for Active	Active
Model Name	MY_Bayesian
Independent Properties	
Calculable Properties	ALogP,Molecular_Weight,Num_H_Donors,Nu...
User Properties	
Cross Validation	True
Folds	5
Learn Options	Validate Models,Remove Uninformative Bins,...
Model Domain Fingerprint	FCFP_2
Additional Properties	
Advanced	
Number of Bins	10

☐ Show Parameter Help

Run Options Cancel Help

# 结果分析

5-Fold Cross-Validation Result									
Model Name	ROC Score	ROC Rating	True Positive	False Negative	False Positive	True Negative	Sensitivity	Specificity	Concordance
MY_Bayesian	0.998	Excellent	171	1	0	174	0.994	1.000	0.997

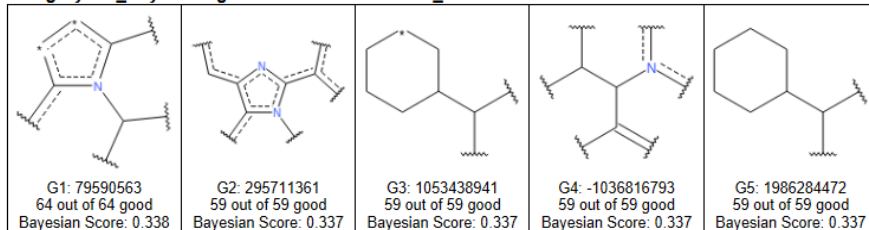


GA

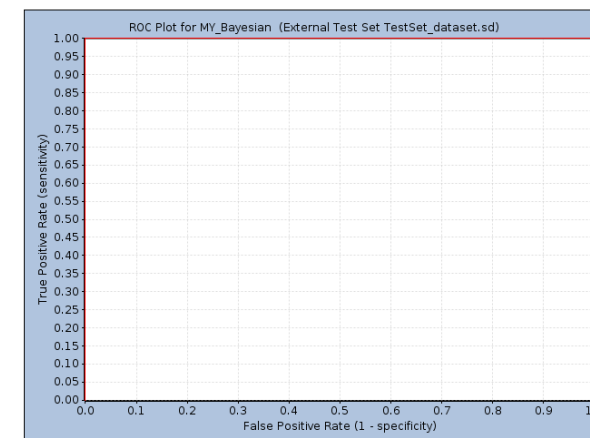
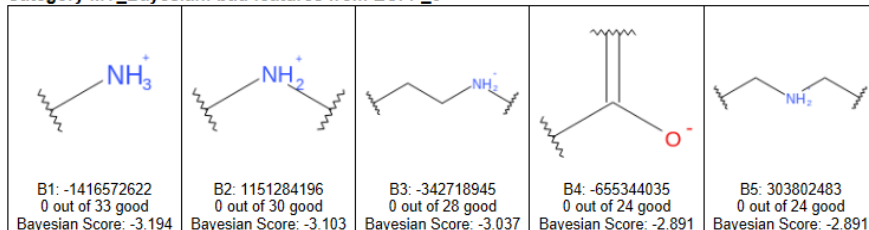
	Index	Name	Visible	Tagged	Visibility Locked	Active	MY_Bayesian	MY_Bayesian#Prediction
40	30	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1	14.2845	true
41	31	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1	11.8302	true
42	32	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1	11.6981	true
43	33	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1	9.0923	true
44	36	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	1	3.7689	true
45	48	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	0	-14.5087	false
46	78	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	0	-25.5528	false
47	45	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	0	-12.6683	false
48	46	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	0	-12.7889	false
49	47	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	0	-14.358	false



Category MY\_Bayesian: good features from ECFP\_6



Category MY\_Bayesian: bad features from ECFP\_6



Validation Result Using External Test Set TestSet_dataset.sd									
Model Name	ROC Score	ROC Rating	True Positive	False Negative	False Positive	True Negative	Sensitivity	Specificity	Concordance
MY_Bayesian	1.000	Excellent	43	1	0	42	0.977	1.000	0.988

# 未知化合物活性预测

Calculate Molecular Properties



## Property Calculation

Basic Arithmetic...

Calculate RMSD...

Calculate Molecular Properties...

Calculate Ligand Efficiency...

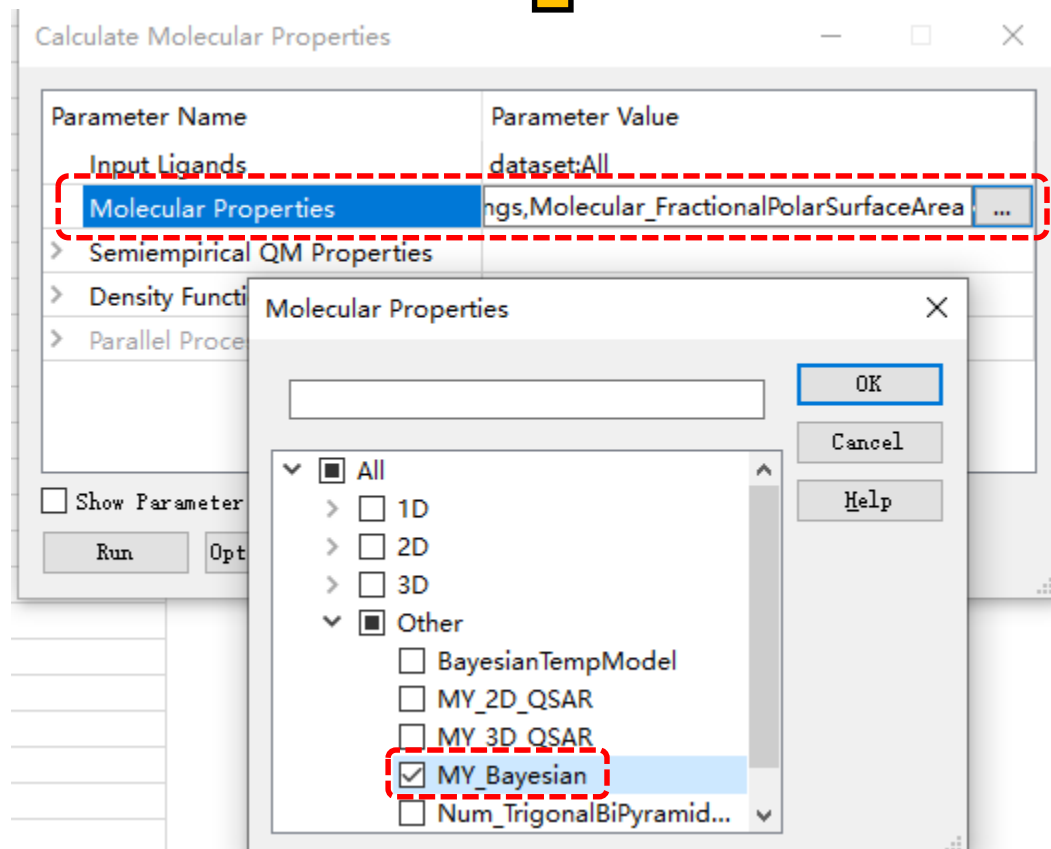


	Index	Name	Visible	Tagged	Visibility Locked	MY_Bayesian	MY_Bayesian#Prediction
1	1	1288666	<input checked="" type="checkbox"/> Yes	<input type="checkbox"/> No	<input type="checkbox"/> No	-3.20069	true
2	2	1320138	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	4.79479	true
3	3	1333937	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	2.05382	true
4	4	2197808	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	9.37912	true
5	5	2260458	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	9.71728	true
6	6	2897692	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	4.33741	true
7	7	2914648	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	5.33722	true
8	8	2954950	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	6.9356	true
9	9	3333	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	5.41247	true
10	10	3567787	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	-8.91515	false



dataset

	Index	Name	Visible	Tagged	Visibility Locked	index
1	1	Molecule	<input checked="" type="checkbox"/> Yes	<input type="checkbox"/> No	<input type="checkbox"/> No	1
2	2	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	2
3	3	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	3
4	4	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	4
5	5	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	5
6	6	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	6
7	7	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	7
8	8	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	8
9	9	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	9
10	10	Molecule	<input type="checkbox"/> No	<input type="checkbox"/> No	<input type="checkbox"/> No	10



## 任务：

- (1) 了解数据集中正负样本处理方法
- (2) 模型构建与结果分析
- (3) 未知化合物活性预测



# 实验报告：

## (1) 实验目的

FXR活性剂判别模型的构建

## (2) 操作流程

正负样本、训练/测试集准备，分子描述符选择，贝叶斯模型构建，模型分析，未知化合物活性预测

## (3) 结果与讨论

模型精度的影响因素、化合物结构分析等