



# Ersilia



## H3D Foundation and Ersilia present

Bringing data science and AI/ML tools to  
infectious disease research

### Event Sponsors



CS&S

Code for  
Science  
& Society





# Session 2: Building a machine learning model

Skills Development workshop

Gemma Turon, @TuronGemma, gemma@ersilia.io

Miquel Duran-Frigola, @mduranfrigola, miquel@ersilia.io

Ersilia Open Source Initiative, @ersiliaio, <https://ersilia.io>

28th September 2022



# Tools

Google Colaboratory



Therapeutics Data Commons



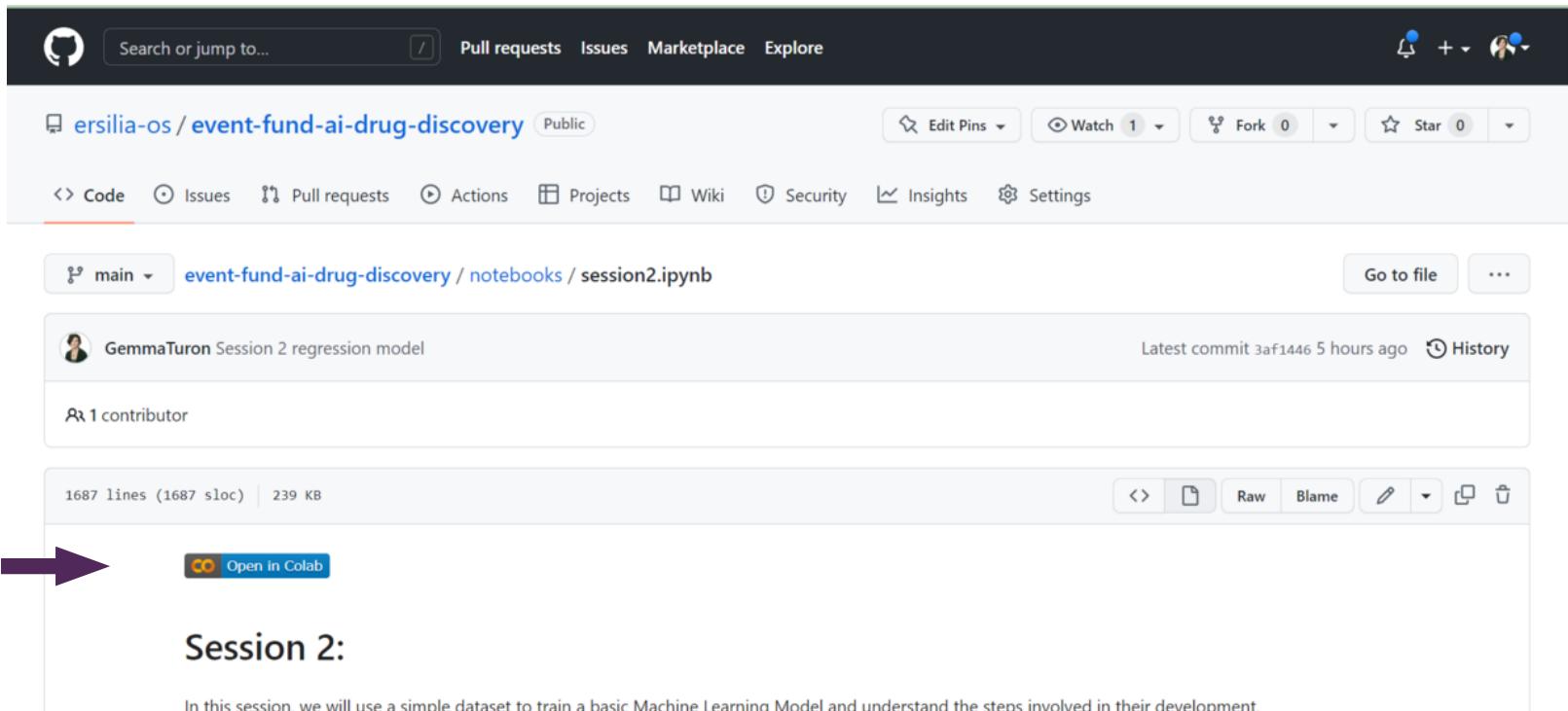
Python Packages





# Getting Started

<https://github.com/ersilia-os/event-fund-ai-drug-discovery>



The screenshot shows a GitHub repository page for 'ersilia-os / event-fund-ai-drug-discovery'. The repository is public and has 1 contributor. A file named 'session2.ipynb' is listed, with a 'Open in Colab' button highlighted by a large purple arrow. The page includes standard GitHub navigation and repository statistics.

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

main event-fund-ai-drug-discovery / notebooks / session2.ipynb Go to file ...

GemmaTuron Session 2 regression model Latest commit 3af1446 5 hours ago History

1 contributor

1687 lines (1687 sloc) | 239 KB

Open in Colab

## Session 2:

In this session, we will use a simple dataset to train a basic Machine Learning Model and understand the steps involved in their development.



# Classification tasks

A classifier identifies the category of a new data point. In the case of bioactivity data, we usually find two categories (binary classifier)

- Active (1)
- Inactive (0)

To learn more about classification, we will use the Therapeutics Data Commons prepared dataset for blockade of hERG activity (Tox)

- Actives (1): block hERG
- Inactives (0): do not block hERG



# Load the data from TDC

## ▶ Download TDC hERG Dataset

- Click on the play button to get the train and tests sets for the TDC hERG data By running this cell, you will get the datasets saved to the drive in data/session2

Show code



```
1 # import the hERG dataset, part of the Toxicity data available at TDC
2 !pip install PyTDC
3
4 from tdc.single_pred import Tox
5 data = Tox(name = "hERG")
6
7 split = data.get_split()
8
9 #we can now separate the compressed dataset in the three sections
10 train = split["train"]
11 validation = split["valid"]
12 test = split["test"]
13
14 train.to_csv("drive/MyDrive/h3d_ersilia_ai_workshop/data/session2/tdc_herg_tr
ain.csv", index=False)
15 validation.to_csv("drive/MyDrive/h3d_ersilia_ai_workshop/data/session2/tdc_he
rg_validation.csv", index=False)
16 test.to_csv("drive/MyDrive/h3d_ersilia_ai_workshop/data/session2/tdc_herg_te
st.csv", index=False)
```



# Load the data from TDC



```
1 #we can check how many molecules we have in each set
2 print(len(train))
3 print(len(valid))
4 print(len(test))
```

458, 66, 131

the train set is used to train the ML model

the validation set is used to improve the model

the test set is used to evaluate the final model performance

In this example we will not focus on optimizing the model parameters, therefore we will only use the train and test splits



# Load the data from TDC



```
1 #we can check how many molecules we have in each set
2 print(len(train))
3 print(len(valid))
4 print(len(test))
```

458, 66, 131

the train set is used to train the ML model

the validation set is used to improve the model

the test set is used to evaluate the final model performance



```
1 train.head()
2 print(len(train[train["Y"]==0]))
3 print(len(train[train["Y"]==1]))
```

	Drug_ID	Drug	Y
0	DEMETHYLASTEMIZOLE	Oc1ccc(CCN2CCC(Nc3nc4ccccc4n3Cc3ccc(F)cc3)CC2)cc1	1.0
1	GBR-12909	Fc1ccc(C(OCC[NH+]2CC[NH+](CCCc3cccc3)CC2)c2cc...	1.0
2	CLOFILUM PHOSPHATE	CCCCCCC[N+](CC)(CC)CCCCc1ccc(Cl)cc1.CCCCCCC[N+...	1.0
3	FLUSPIRILENE	O=C1NCN(c2ccccc2)C12CC[NH+](CCCC(c1ccc(F)cc1)c...	1.0
4	VANOXERINE HYDROCHLORIDE	Fc1ccc(C(OCCN2CCN(CCCc3cccc3)CC2)c2ccc(F)cc2)cc1	1.0



# Molecular Representation



```
1 !pip install rdkit #cheminformatics package
2 from rdkit import Chem
3 from rdkit.Chem import Draw
4
5 smiles = train["Drug"][449:] # select a few smiles for representation
6 #convert the smiles to RdKit Molecules (not readable by humans)
7 mols = [Chem.MolFromSmiles(smi) for smi in smiles]
```

<rdkit.Chem.rdchem.Mol at 0x7fa16c9435d0>



```
1 from rdkit.Chem import Draw
2 Draw.MolsToGridImage(mols)
```





# Molecular Featurization

Featurization is the conversion of SMILES into vectors or images that can be passed to the ML model

```
● ● ●  
1 !pip install signaturizer  
2 from signaturizer import Signaturizer  
3 sign = Signaturizer("GLOBAL")  
4 X_train = sign.predict(train["Drug"]).signature  
5 X_test = sign.predict(test["Drug"]).signature
```

```
● ● ●
```

```
1 X_train[0]
```

```
array([-0.09308645,  0.09305048, -0.09294817, ..., -0.05613168,  
       -0.1379203 ,  0.05799835], dtype=float32)
```



## Data we have so far:

Original datasets: table with SMILES and its associated hERG blocking activity as 0 or 1

- train, test

X datasets: vector-like representation of the molecules calculated using the Chemical Checker --> model input

- X\_train, X\_test

Y\_datasets: list of 0 and 1 results that correspond to the list of molecules in the train and test sets

- Y\_train, Y\_test



# Supervised Machine Learning

In summary, we now have two datasets (train, test):

- X (input): molecular signature
- Y (output): activity against hERG, 1 or 0

We will use the Random Forest Classifier algorithm implemented in scikit-learn.

```
● ● ●  
1 !pip install sklearn  
2 from sklearn.ensemble import RandomForestClassifier  
3 clf = RandomForestClassifier()  
4 clf.fit(X_train, Y_train)
```



# Classification outputs

Once the model is fitted, we can predict the category of each molecule in the validation and test sets. A classifier outputs two numbers per each prediction:

- Probability of 0
- Probability of 1



```
1 y_pred = clf.predict_proba(x_test)
2 y_pred[:10]
```



```
1 y_pred[:, 1]
```

```
array([[0.31, 0.69],
       [0.37, 0.63],
       [0.69, 0.31],
       [0.14, 0.86],
       [0.58, 0.42],
       [0.84, 0.16],
       [0.163, 0.837],
       [0.257, 0.743],
       [0.111, 0.889],
       [0.303, 0.697]])
```



# Classification outputs

In a classification, the output is a probability. We need to define a cut-off or threshold to transform the results into a binary output (0 or 1) again. The threshold is typically set at 0.5 by default.

Proba 0	Proba 1	Cut-off 0.5	Cut-off 0.7	Cut-off 0.3
0.39	0.61			
0.3	0.7			
0.69	0.31			
0.21	0.79			



# Classification outputs

In a classification, the output is a probability. We need to define a cut-off or threshold to transform the results into a binary output (0 or 1) again. The threshold is typically set at 0.5 by default.

Proba 0	Proba 1	Cut-off 0.5	Cut-off 0.7	Cut-off 0.3
0.39	0.61	1	0	1
0.3	0.7			
0.69	0.31			
0.21	0.79			



# Classification outputs

In a classification, the output is a probability. We need to define a cut-off or threshold to transform the results into a binary output (0 or 1) again. The threshold is typically set at 0.5 by default.

Proba 0	Proba 1	Cut-off 0.5	Cut-off 0.7	Cut-off 0.3
0.39	0.61	1	0	1
0.3	0.7	1	1	1
0.69	0.31			
0.21	0.79			



# Classification outputs

In a classification, the output is a probability. We need to define a cut-off or threshold to transform the results into a binary output (0 or 1) again. The threshold is typically set at 0.5 by default.

Proba 0	Proba 1	Cut-off 0.5	Cut-off 0.7	Cut-off 0.3
0.39	0.61	1	0	1
0.3	0.7	1	1	1
0.69	0.31	0	0	1
0.21	0.79			



# Classification outputs

In a classification, the output is a probability. We need to define a cut-off or threshold to transform the results into a binary output (0 or 1) again. The threshold is typically set at 0.5 by default.

Proba 0	Proba 1	Cut-off 0.5	Cut-off 0.7	Cut-off 0.3
0.39	0.61	1	0	1
0.3	0.7	1	1	1
0.69	0.31	0	0	1
0.21	0.79	1	1	1



# Model Evaluation

## Confusion Matrix

		Inactives	Actives
True	Inactives	True Negative	False Positive
	Actives	False Negative	True Positive
Predicted		Inactives	Actives

Precision: how many positives are actually positive

$$\frac{\text{TP}}{\text{TP}+\text{FP}}$$

Recall: how many positives are we able to identify

$$\frac{\text{TP}}{\text{TP}+\text{FN}}$$

Precision

Recall





# Model Evaluation

## Confusion Matrix

	Inactives	True Negative	False Positive	
True	Actives	False Negative	True Positive	
	Inactives	Actives		
Predicted				

Higher, more restrictive threshold in proba1

Precision

Recall



Lower threshold in proba1

Recall

Precision





# Model Evaluation

## Confusion Matrix

		Inactives	Actives
True	Inactives	True Negative	False Positive
	Actives	False Negative	True Positive
		Inactives	Actives
Predicted			

In which scenarios we might want a high precision?

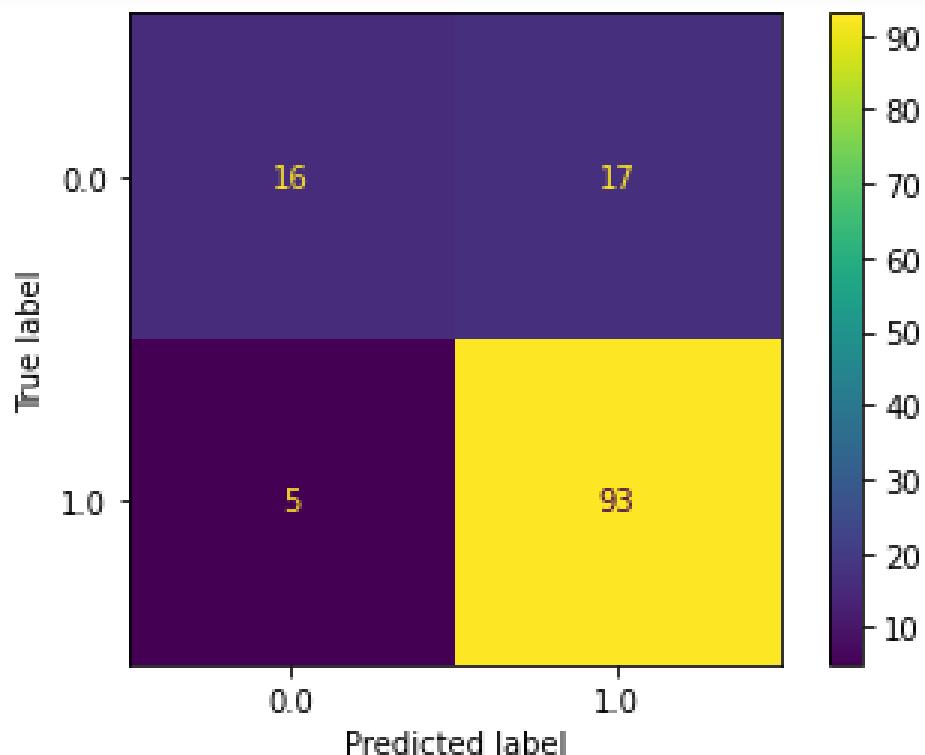
In which scenarios we might want a high recall?



# Model Evaluation

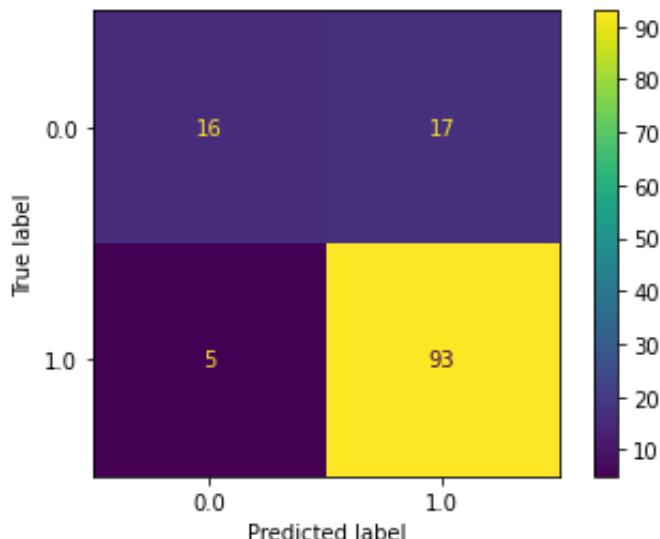


```
1 from sklearn.metrics import ConfusionMatrixDisplay as cdm
2 cdm.from_estimator(clf, X_test, Y_test)
```





# Model Evaluation



Precision: how many positives are actually positive

Recall: how many positives are we able to identify

$$\frac{\text{TP}}{\text{TP}+\text{FP}} = \frac{93}{93+17}$$

$$\frac{\text{TP}}{\text{TP}+\text{FN}} = \frac{93}{93+5}$$



# Model Evaluation



```
1 from sklearn.metrics import precision_score, recall_score
2 #precision and recall scores need a specific threshold
3
4 y_pred_bin = []
5 for x in y_pred[:,1]:
6     if x > 0.5:
7         y_pred_bin += [1]
8     if x <= 0.5:
9         y_pred_bin += [0]
10
11 precision = precision_score(Y_test, y_pred_bin)
12 recall = recall_score(Y_test, y_pred_bin)
13 print(precision, recall)
```

0.84545454545455 0.9489795918367347



# Data we have so far:

Original datasets: train, test

X datasets: X\_train, X\_test

Y\_datasets: Y\_train, Y\_test

Predictions:

- y\_pred: output of the classifier, expressing the probability that a molecule is inactive (0) or active (1).
- y\_pred\_bin: binarized activity based on the predicted probability, the cut-off is set by the researcher.



# Model Evaluation

ROC Curve: performance of the model at all classification thresholds (0-1)

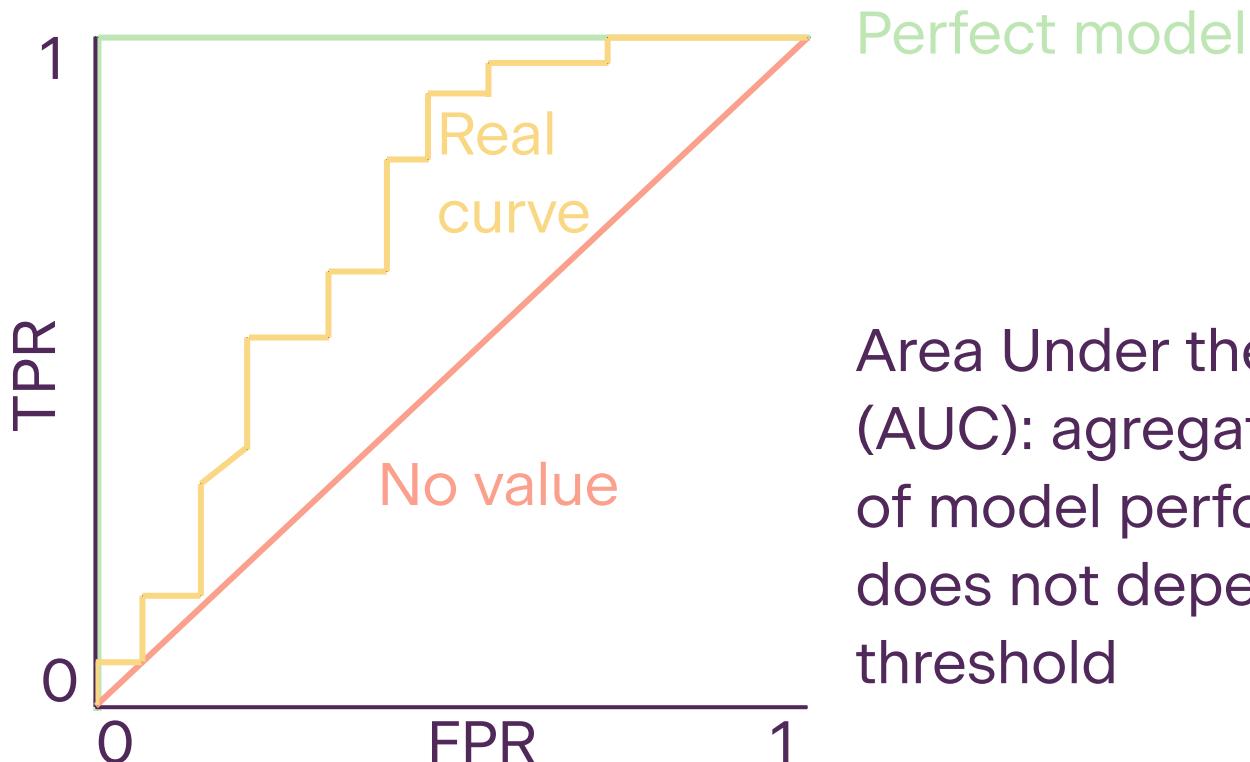
		True	
	Inactives	True Negative	False Positive
True			
Actives		False Negative	True Positive
	Inactives	Actives	
		Predicted	

True Positive Rate (sensitivity or recall): proportion of correctly predicted positives of all positive observations ( $TP/(TP+FN)$ )

False Positive Rate (100-sensitivity): proportion of incorrectly predicted positives of all negative observations ( $FP/(TN+FP)$ )

# Model Evaluation

ROC Curve: performance of the model at all classification thresholds (0-1)



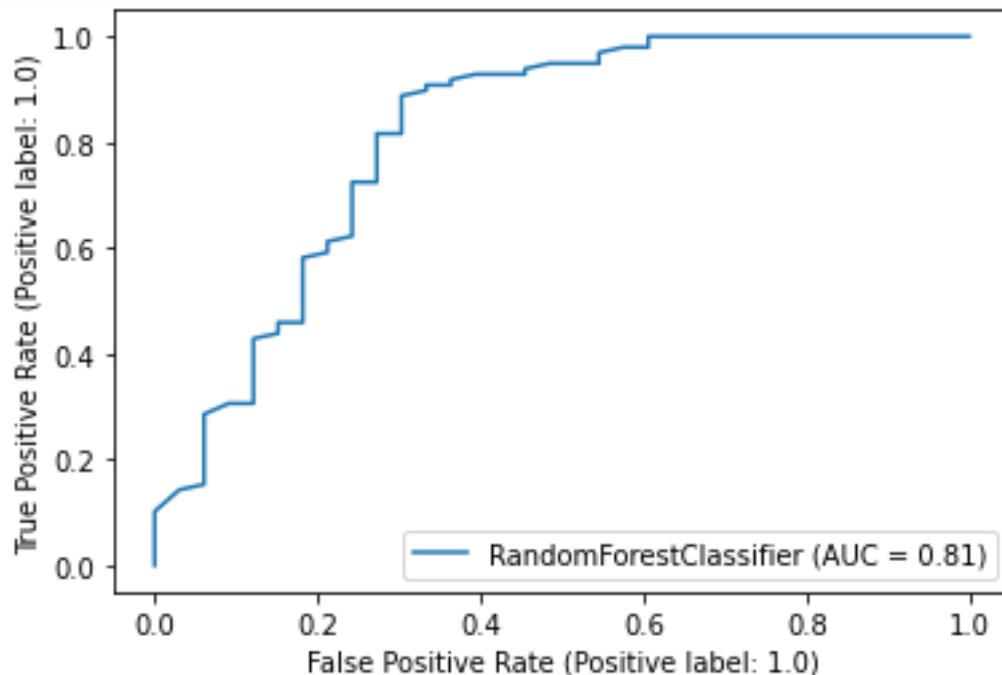
Area Under the Curve (AUC): aggregate measure of model performance. It does not depend on the threshold



# Model Evaluation



```
1 from sklearn.metrics import RocCurveDisplay as rdc
2 rdc.from_estimator(clf, X_test, Y_test)
```

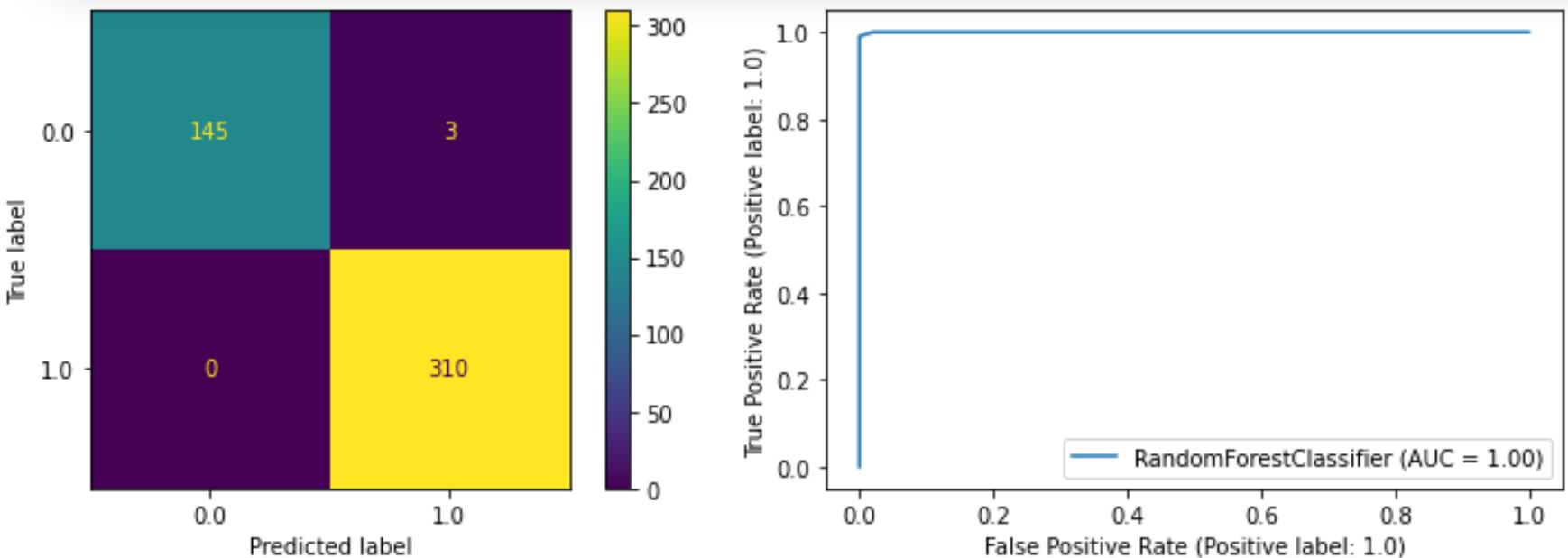




# Model Evaluation



```
1 from sklearn.metrics import ConfusionMatrixDisplay as cdm
2 cdm.from_estimator(clf, X_train, Y_train)
3
4 from sklearn.metrics import RocCurveDisplay as rdc
5 rdc.from_estimator(clf, X_train, Y_train)
```





# Classification tasks: summary

- We have used a binary classification for hERG activity and trained a random forest classifier
- The molecules were featurized using Chemical Checker signatures, which combine structural and bioactivity data of the molecules
- The model performance was good (AUC above 0.8) on both the validation and test sets
- How could we try to improve the model?



# Regression tasks

A regressor predicts continuous numerical values, for example:

- IC50
- % of inhibition

To learn more about regression, we will use the Therapeutics Data Commons Acute Toxicity:

- Lethal Dose 50 (LD50): amount of drug (mg/kg) that kills 50% of the tested population



# Load the data from TDC

## ▶ Download TDC LD50 Dataset

● Click on the play button to get the train and tests sets for the TDC hERG data By running this cell, you will get the datasets saved to the drive in data/session2

[Show code](#)



```
1
2 #we can check how many molecules we have in each set
3 print(len(train))
4 print(len(valid))
5 print(len(test))
```

Train set: 5170 molecules

Validation set:  
738 molecules

Test set: 1477  
molecules



# Molecular Representation

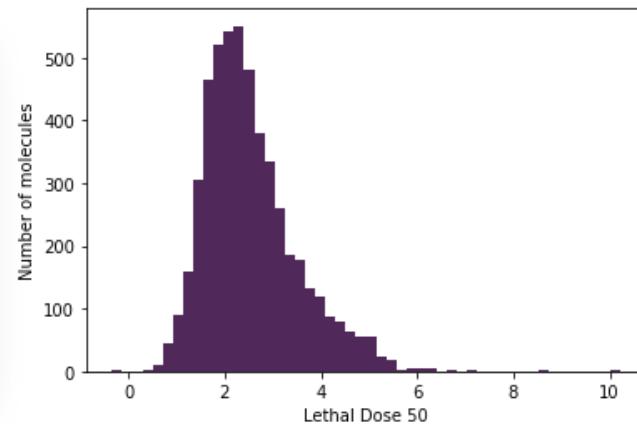


```
1 train.head()
```

	Drug_ID	Drug	Y
0	Methane, tribromo-	BrC(Br)Br	2.343
1	Bromoethene (9Cl)	C=CBr	2.330
2	1,1'-Biphenyl, hexabromo-	Brc1ccc(-c2ccc(Br)c(Br)c2Br)c(Br)c1Br	1.465
3	Isothiocyanic acid, p-bromophenyl ester	S=C=Nc1ccc(Br)cc1	2.729
4	Benzene, bromo-	Brc1ccccc1	1.765



```
1 import matplotlib.pyplot as plt
2
3 plt.hist(train["Y"], bins=50, color
4         = "#50285a")
5 plt.xlabel("Lethal Dose 50")
6 plt.ylabel("Number of molecules")
```





# Molecular Featurization

## Chemical Checker Signatures



```
1 from signaturizer import Signaturizer
2 sign = Signaturizer("GLOBAL")
3 X_train = sign.predict(train["Drug"]).signature
4 X_test = sign.predict(test["Drug"]).signature
```

```
array([-0.09308645,  0.09305048, -0.09294817, ..., -0.05613168,
       -0.1379203 ,  0.05799835], dtype=float32)
```



```
1 Y_train = list(train["Y"])
2 Y_test = list(test["Y"])
```

2.343



# Supervised Machine Learning

In summary, we now have three datasets (train, validation, test):

- X (input): molecular signature
- Y (output): LD50 values

We will use a simple linear regression from sci-kit learn to train a model



```
1 from sklearn.linear_model import LinearRegression
2 reg = LinearRegression()
3 reg.fit(X_train, Y_train)
```



# Regression outputs

The regression gives back the predicted value, in our case, LD50.

```
● ● ●  
1 y_pred = reg.predict(x_test)  
2  
3 #if we check, the prediction is simply a list of LD50  
4 y_pred[:4]
```

```
array([3.3221316, 2.0938659, 2.874437 , 4.17527 ], dtype=float32)
```

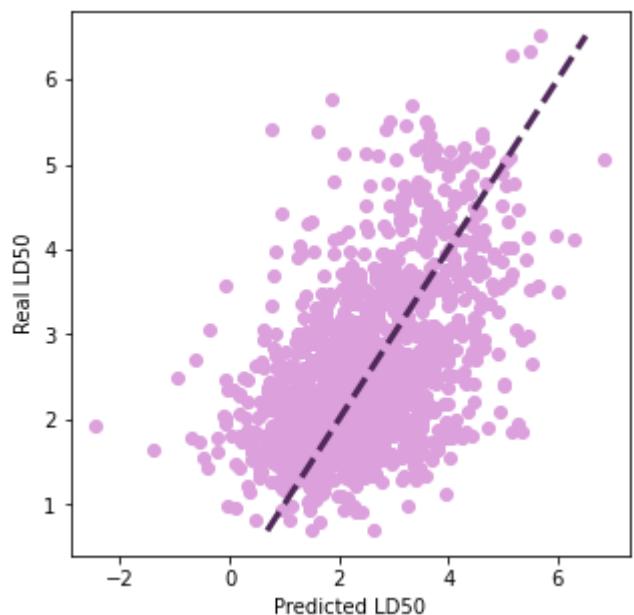


# Regression model evaluation

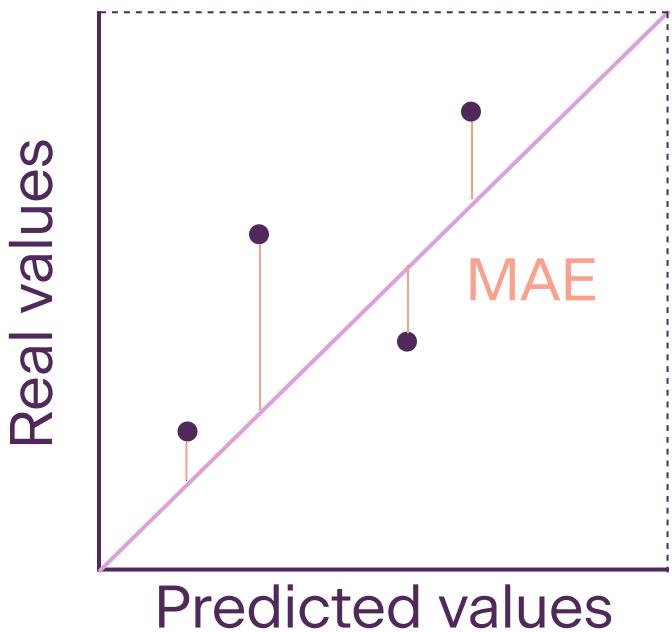
The easiest is to plot a scatter distribution of the real versus the predicted values.

We use the Matplotlib package for plotting functions

```
● ● ●  
1 fig, ax = plt.subplots(1,1, figsize=(5,5))  
2 ax.scatter(x=y_pred, y=Y_test,  
3 color="#dca0dc")  
4 ax.set_xlabel("Predicted LD50")  
5 ax.set_ylabel("Real LD50")
```



# Regression Model Evaluation



- Mean Absolute Error (MAE): average of the absolute difference between actual and predicted values
- Mean Squared Error (MSE): average of the squared difference
- R-square: proportion of the variance explained by the regression model. It is scale free



# Regression model evaluation



```
1 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
2
3 mae = mean_absolute_error(Y_valid, valid_pred)
4 mse = mean_squared_error(Y_valid, valid_pred)
5 r2 = r2_score(Y_valid, valid_pred)
6
7 print("MAE: {}".format(mae))
8 print("MSE: {}".format(mse))
9 print("R2: {}".format(r2))
```

MAE: 0.8202228248986185

MSE: 1.1181765214543737

R2: -0.25158525929946496

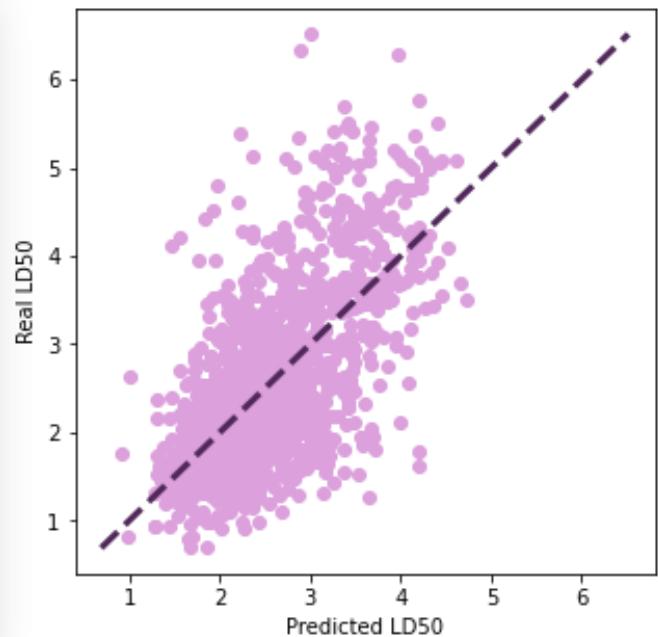
What do we think of these results?



# Improving regression models

We can test other models for regression to improve the performance on our data

```
● ● ●  
1 from sklearn.ensemble import  
    RandomForestRegressor  
2  
3 reg = RandomForestRegressor(n_estimators=10)  
    #we use only 10 trees for speed time  
4 reg.fit(X_train, Y_train)  
5  
6 test_pred = reg.predict(X_test)  
7  
8 fig, ax = plt.subplots(1,1, figsize=(5,5))  
9 ax.scatter(x=y_pred, y=Y_test,  
    color="#dca0dc")  
10 ax.set_xlabel("Predicted LD50")  
11 ax.set_ylabel("Real LD50")
```





# Improving regression models

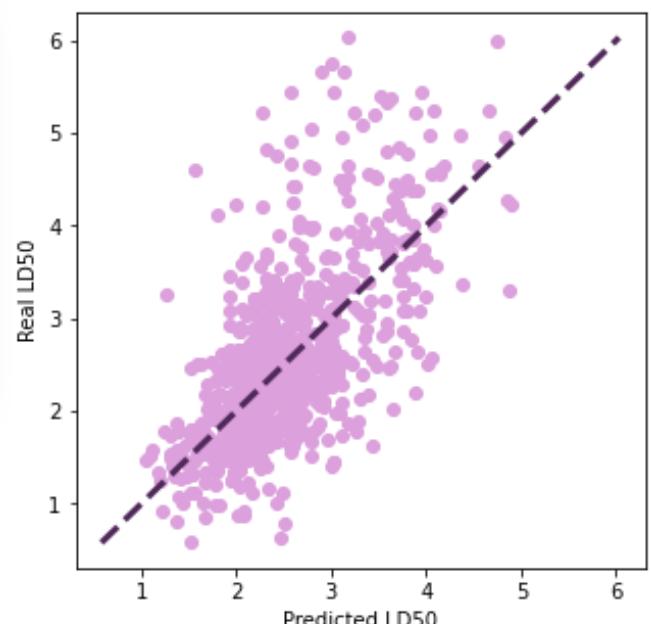
We can test other models for regression to improve the performance on our data

```
1 mae = mean_absolute_error(Y_valid, valid_pred)
2 mse = mean_squared_error(Y_valid, valid_pred)
3 r2 = r2_score(Y_valid, valid_pred)
4
5 print("MAE: {}".format(mae))
6 print("MSE: {}".format(mse))
7 print("R2: {}".format(r2))
```

MAE: 0.5450362715297566

MSE: 0.5291154421053693

R2: 0.40775622175874715





# Conclusions

- We have used the Therapeutics Data Commons datasets to build simple classification and regression QSAR models:
  - Random Forest Classifier
  - Linear Regression
  - Random Forest Regressor
- We have evaluated model performance using the validation and test sets, and the basic metrics for:
  - Classification: precision, recall and AUC
  - Regression: MAE, MSE, R2