

# KunitzProject\_Tornisiello

June 11, 2020

- 1 This notebook contains the python script used to handle the datasets, compute the MCC, chose the thresholds, compute the accuracy, generate confusion matrixes, ROC curves and strip-plot

```
[1]: #creation of a pandas series that contains the UniProt ID and the class
      ↪ (1=postive, 0=negative)
import pandas as pd
import numpy as np
from Bio import SeqIO
from Bio import SearchIO

#domain series creation
domain = {}
for record in SeqIO.parse("./negatives_kunitz2.fasta", "fasta"):
    ID_n = record.id.split("|")[1]
    domain[ID_n] = 0
for record in SeqIO.parse("./positives_clean_kunitz2.fasta", "fasta"):
    ID_p = record.id.split("|")[1]
    domain[ID_p] = 1
domain = pd.Series(domain, dtype='category')
domain
```

```
[1]: Q4R8P0      0
     P03949      0
     Q9NPB9      0
     P31937      0
     Q9FFR3      0
     ..
     PODJ49      1
     Q90W98      1
     PODJ68      1
     PODJ65      1
     PODJ77      1
     Length: 562240, dtype: category
     Categories (2, int64): [0, 1]
```

```
[2]: #creation of a pandas series containing the UniProtID and the full sequence
      ↪E-value
hmm_results = SearchIO.read("./hmmsearch_output_kunitz.txt", "hmmer3-tab")
evaluate = {}
for hit in hmm_results:
    ID = hit.id.split("|")[1]
    evaluate[ID] = hit.evalue
evaluate = pd.Series(evaluate)
evaluate
```

```
[2]: Q868Z9      1.200000e-194
      O76840      2.400000e-179
      Q02445      8.500000e-68
      P84875      6.000000e-67
      O54819      4.300000e-66
      ...
      W4VS46      1.000000e+00
      W5U5X5      1.000000e+00
      X5CFH4      1.000000e+00
      X5CWH9      1.000000e+00
      X5IWT5      1.000000e+00
      Length: 269842, dtype: float64
```

```
[3]: #creation of a unique pandas dataframe containing ID, class and full seq E-value
dataframe = pd.DataFrame({'EVALUE': evaluate, 'CLASS': domain})
dataframe['EVALUE'].fillna(10, inplace = True)
```

```
[4]: #split of the dataframe in train and test subsets
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
import seaborn as sns
#full seq e-value=A
test_size = 0.5
seed = 42
X_trainA, X_testA, Y_trainA, Y_testA = train_test_split(dataframe.EVALUE,
      ↪dataframe.CLASS, test_size=test_size, random_state=seed)
X_trainA_2D = X_trainA.values.reshape(-1, 1)
X_testA_2D = X_testA.values.reshape(-1, 1)
```

```
[5]: #computation of MCC for each E-value (full sequence) threshold (ranging from
      ↪1e-20 to 1)
from sklearn.metrics import matthews_corrcoef

mcc_list = []
for i in range(-20, 1, 1):
    Y_pred = X_trainA.apply(lambda x: 1 if x<10**(i) else 0)
    mcc_i = matthews_corrcoef(Y_trainA, Y_pred)
```

```

    mcc_list.append((i, mcc_i))
    Y_pred = Y_pred.iloc[0:0]
print(mcc_list)

```

```

[(-20, 0.9196032064517125), (-19, 0.9591604070016291), (-18,
0.9739480125404418), (-17, 0.9768787384319279), (-16, 0.9827140355715218), (-15,
0.9856187616254077), (-14, 0.9885149739771153), (-13, 0.9914027472212671), (-12,
0.9942821548701704), (-11, 0.9942821548701704), (-10, 0.9942821548701704), (-9,
0.9971532693756712), (-8, 0.994330692767468), (-7, 0.994330692767468), (-6,
0.9915319304689825), (-5, 0.9725802156114983), (-4, 0.942472976225046), (-3,
0.7161656597018671), (-2, 0.2941020423004492), (-1, 0.0667530101845895), (0,
0.026873608326550186)]

```

```

[10]: #best th=10e-9 full sequence evalve
#application on test subset and computation of the stripplot
Y_predA = X_testA.apply(lambda x: 1 if x<10e-9 else 0)
mismatch = Y_predA.loc[Y_testA != Y_predA]
false_pos = mismatch.loc[mismatch == 1].index
false_neg = mismatch.loc[mismatch == 0].index
print('false positives are:', false_pos)
print('false negatives are:', false_neg)
striplot = sns.stripplot(y = X_testA, x = Y_testA)
striplot.set(yscale='log', ylim=(10e-27,40))
striplot.axhline(10e-9)

```

```

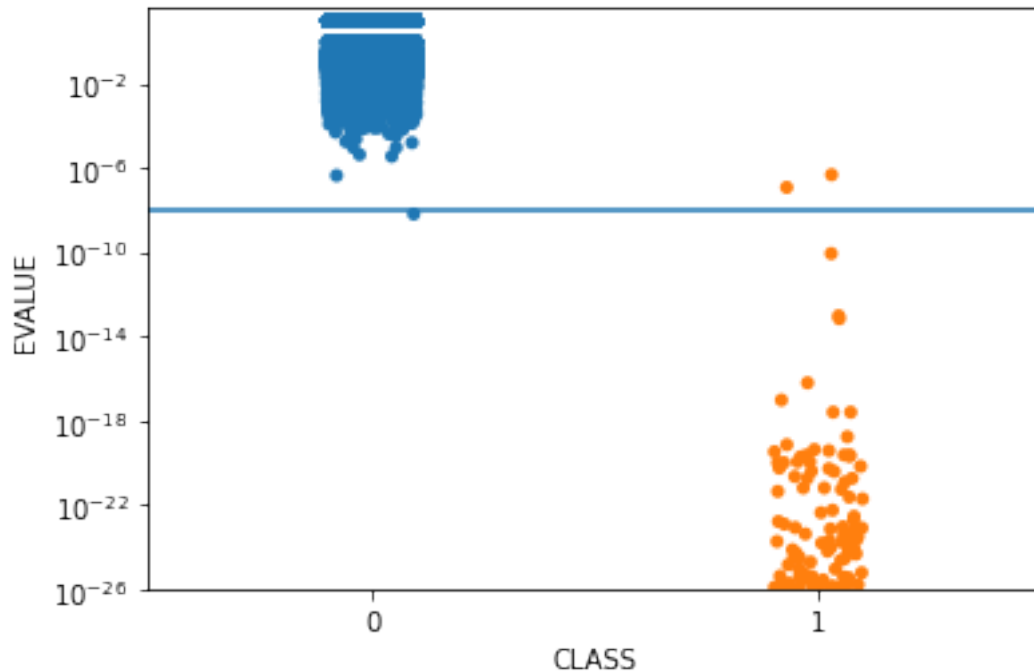
false positives are: Index(['P84555'], dtype='object')
false negatives are: Index(['D3GGZ8', '062247'], dtype='object')

```

```

[10]: <matplotlib.lines.Line2D at 0x7fcc8c559f40>

```



```
[11]: import sklearn
sklearn.metrics.accuracy_score(Y_testA, Y_predA)
```

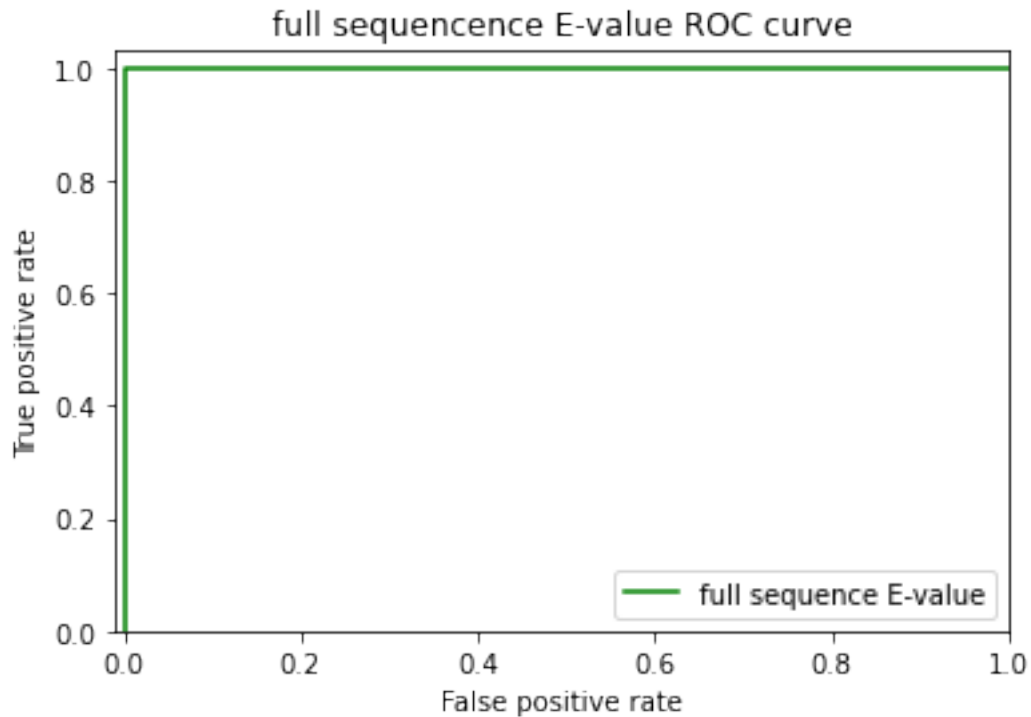
```
[11]: 0.999989328400683
```

```
[12]: #computation of confusion matrix full sequence evalve
from sklearn.metrics import confusion_matrix
confusion_matrix(Y_testA, Y_predA)
```

```
[12]: array([[280948,    1],
           [    2,   169]])
```

```
[13]: #computation of ROC curve and AUC for full sequence e-value
import sklearn
from sklearn import metrics
import matplotlib.pyplot as plt
Y_score = [- values for values in X_testA]
fprA, tprA, thA = sklearn.metrics.roc_curve(Y_testA, Y_score)
aucA = metrics.auc(fprA, tprA)
#plot ROC curve
plt.title('full sequence E-value ROC curve')
plt.plot(fprA, tprA, color='green', label='full sequence E-value')
plt.legend(loc='lower right')
plt.xlim([-0.01, 1])
plt.ylim([0, 1.03])
```

```
plt.ylabel('True positive rate')
plt.xlabel('False positive rate')
plt.show()
aucA
```



```
[13]: 0.9999999375550024
```

### 1.1 Repeating the same process for best domain E-value

```
[14]: #evaluate-best domain pandas series
hmm_results = SearchIO.read("./hmmsearch_output_kunitz.txt", "hmmer3-tab")
evaluate = {}
for hit in hmm_results:
    ID = hit.id.split("|")[1]
    evaluate[ID] = hit[0].evaluate
evaluate = pd.Series(evaluate)
evaluate
```

```
[14]: Q868Z9      1.400000e-22
      O76840      5.100000e-23
      Q02445      2.300000e-26
      P84875      1.400000e-25
```

```

054819    1.600000e-25
...
W4VS46    2.600000e-01
W5U5X5    1.000000e+00
X5CFH4    1.000000e+00
X5CWH9    1.000000e+00
X5IWT5    1.000000e+00
Length: 269842, dtype: float64

```

```

[15]: #creation of the dataframe
dataframe = pd.DataFrame({'EVALUE': evaluate, 'CLASS': domain})
dataframe['EVALUE'].fillna(10, inplace = True)

```

```

[16]: #splitting in train and test subsets
#best domain e-value=B
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import plot_confusion_matrix
import seaborn as sns

test_size = 0.5
seed = 42
X_trainB, X_testB, Y_trainB, Y_testB = train_test_split(dataframe.EVALUE,
↳dataframe.CLASS, test_size=test_size, random_state=seed)
X_trainB_2D = X_trainB.values.reshape(-1, 1)
X_testB_2D = X_testB.values.reshape(-1, 1)

```

```

[17]: #computation of MCC for best domain evaluate for each treshold (ranging from
↳1e-20 to 1)
mcc_list = []
for i in range(-20, 1, 1):
    Y_predB = X_trainB.apply(lambda x: 1 if x<10**(i) else 0)
    mcc_i = matthews_corrcoef(Y_trainB, Y_predB)
    mcc_list.append((i, mcc_i))
    Y_predB = Y_predB.iloc[0:0]
print(mcc_list)

```

```

[(-20, 0.9164897961982789), (-19, 0.9591604070016291), (-18,
0.9739480125404418), (-17, 0.9768787384319279), (-16, 0.9827140355715218), (-15,
0.9827140355715218), (-14, 0.9885149739771153), (-13, 0.9914027472212671), (-12,
0.9942821548701704), (-11, 0.9942821548701704), (-10, 0.9942821548701704), (-9,
0.9971532693756712), (-8, 0.9971532693756712), (-7, 0.994330692767468), (-6,
0.9915319304689825), (-5, 0.9860045230665895), (-4, 0.952196228309818), (-3,
0.8005011159991144), (-2, 0.4072229500937661), (-1, 0.12766159072532623), (0,
0.03989974682744708)]

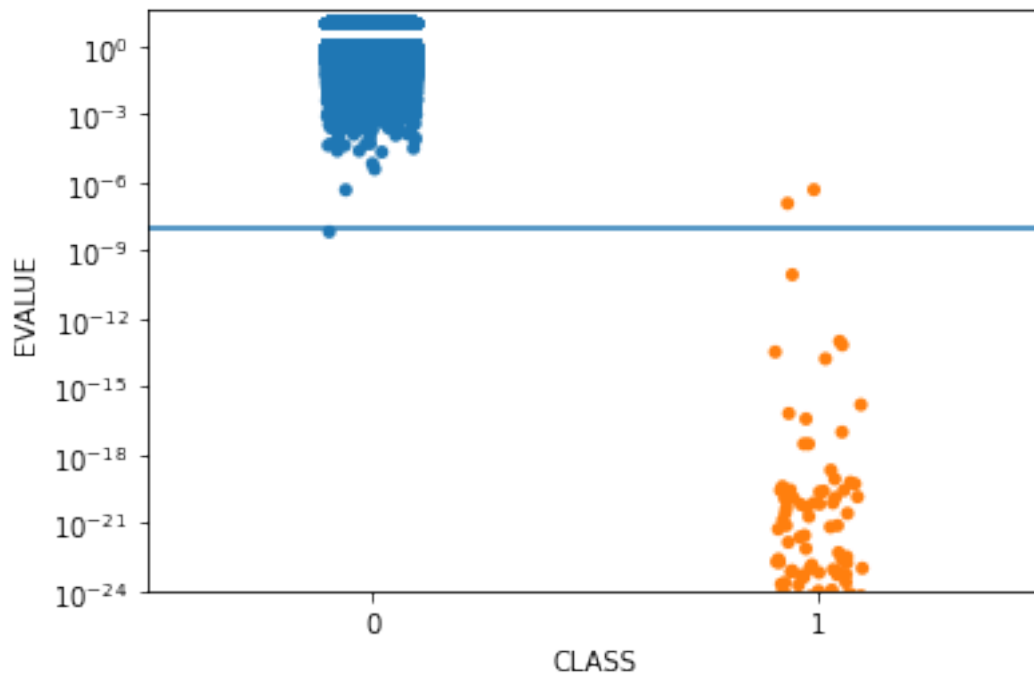
```

```
[26]: #best th= 10e-9 best domain evaluate
#on test
Y_predB = X_testB.apply(lambda x: 1 if x<10e-9 else 0)
mismatch = Y_predB.loc[Y_testB != Y_predB]
false_pos = mismatch.loc[mismatch == 1].index
false_neg = mismatch.loc[mismatch == 0].index
print('false positives are:', false_pos)
print('false negatives are:', false_neg)
striplot = sns.striplot(y = X_testB, x = Y_testB)
striplot.set(yscale='log', ylim=(10e-25,40))
striplot.axhline(10e-9)
```

false positives are: Index(['P84555'], dtype='object')

false negatives are: Index(['D3GGZ8', '062247'], dtype='object')

```
[26]: <matplotlib.lines.Line2D at 0x7fcc9500b250>
```



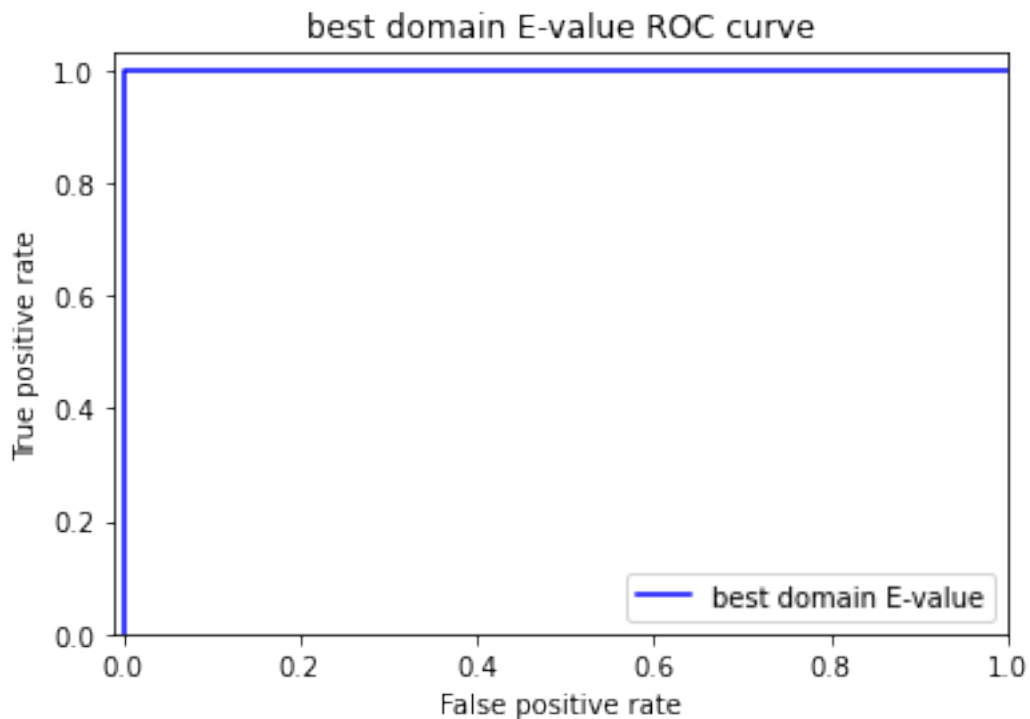
```
[27]: sklearn.metrics.accuracy_score(Y_testB, Y_predB)
```

```
[27]: 0.999989328400683
```

```
[28]: #confusion matrix best domain
from sklearn.metrics import confusion_matrix
confusion_matrix(Y_testB, Y_predB)
```

```
[28]: array([[280948,    1],
           [    2,   169]])
```

```
[29]: #computation of ROC curve and ACU for best domain e-value
import sklearn
from sklearn import metrics
import matplotlib.pyplot as plt
Y_score = [- values for values in X_testB]
fprB, tprB, thB = sklearn.metrics.roc_curve(Y_testB, Y_score)
aucB = metrics.auc(fprB, tprB)
#plot ROC curve
plt.title('best domain E-value ROC curve')
plt.plot(fprB, tprB, color='blue', label='best domain E-value')
plt.legend(loc='lower right')
plt.xlim([-0.01, 1])
plt.ylim([0, 1.03])
plt.ylabel('True positive rate')
plt.xlabel('False positive rate')
plt.show()
aucB
```



```
[29]: 0.9999999375550024
```



## 1.2 Switching train and test datasets

```
[30]: X_trainA2 = X_testA
      Y_trainA2 = Y_testA
      X_testA2 = X_trainA
      Y_testA2 = Y_trainA
      X_trainA2_2D = X_testA_2D
      X_testA2_2D = X_trainA_2D
```

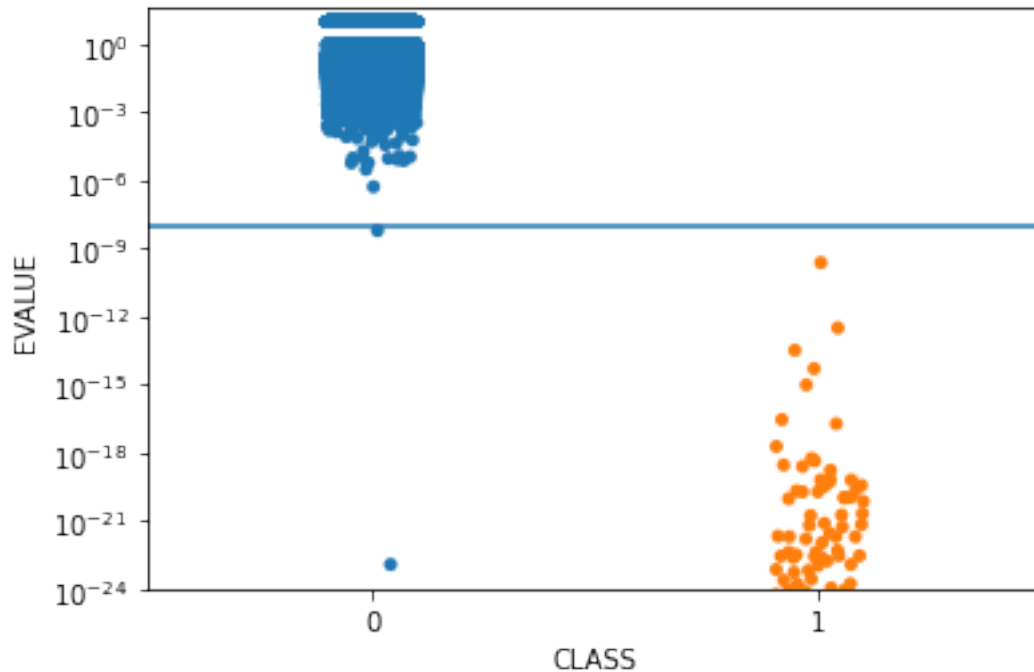
```
[31]: #computation of MCC for each E-value (full sequence) threshold (ranging from
      →1e-20 to 1)
      mcc_list = []
      for i in range(-20, 1, 1):
          Y_predA2 = X_trainA2.apply(lambda x: 1 if x<10**(i) else 0)
          mcc_i = matthews_corrcoef(Y_trainA2, Y_predA2)
          mcc_list.append((i, mcc_i))
          Y_predA2 = Y_predA2.iloc[0:0]
      print(mcc_list)
```

```
[(-20, 0.9302828750256974), (-19, 0.9703025079222909), (-18,
0.9733129372179947), (-17, 0.9822889973386646), (-16, 0.985262887647878), (-15,
0.985262887647878), (-14, 0.985262887647878), (-13, 0.9911839636027662), (-12,
0.9911839636027662), (-11, 0.9911839636027662), (-10, 0.9941313083026798), (-9,
0.9941313083026798), (-8, 0.9912012743277198), (-7, 0.9912012743277198), (-6,
0.9941993109664184), (-5, 0.9828942276077673), (-4, 0.9388136871967353), (-3,
0.7131827817045778), (-2, 0.28295364896244113), (-1, 0.06567514060545633), (0,
0.02651933480552863)]
```

```
[32]: #th= 10e-9 full seq eval
      #on test
      Y_predA2 = X_testA2.apply(lambda x: 1 if x<10e-9 else 0)
      mismatch = Y_predA2.loc[Y_testA2 != Y_predA2]
      false_pos = mismatch.loc[mismatch == 1].index
      false_neg = mismatch.loc[mismatch == 0].index
      print('false positives are:', false_pos)
      print('false negatives are:', false_neg)
      striplot = sns.striplot(y = X_testA2, x = Y_testA2)
      striplot.set(yscale='log', ylim=(10e-25,40))
      striplot.axhline(10e-9)
```

```
false positives are: Index(['P56409', 'G3LH89'], dtype='object')
false negatives are: Index([], dtype='object')
```

```
[32]: <matplotlib.lines.Line2D at 0x7fcc97c95790>
```



```
[33]: sklearn.metrics.accuracy_score(Y_testA2, Y_predA2)
```

```
[33]: 0.9999928856004553
```

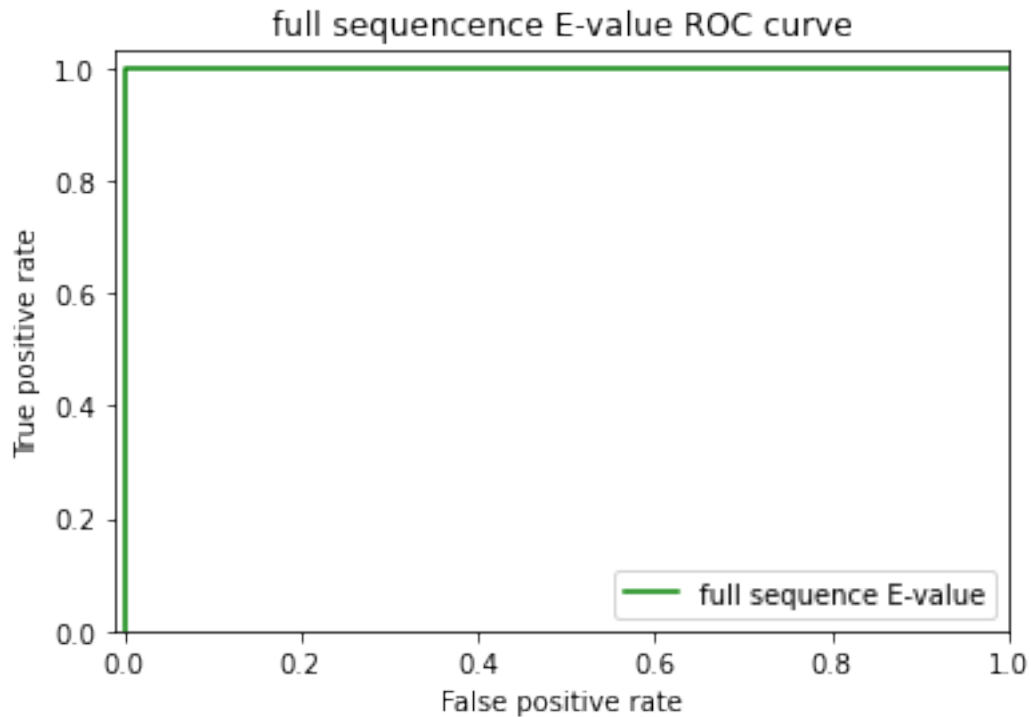
```
[34]: #confusion matrix full seq
from sklearn.metrics import confusion_matrix
confusion_matrix(Y_testA2, Y_predA2)
```

```
[34]: array([[280943,    2],
           [    0,   175]])
```

```
[66]: #computation of ROC curve and ACU for full seq E-value
import sklearn
from sklearn import metrics
import matplotlib.pyplot as plt

Y_score = [- values for values in X_testA2]
fprA2, tprA2, thA2 = sklearn.metrics.roc_curve(Y_testA2, Y_score)
aucA2 = metrics.auc(fprA2, tprA2)
#plot ROC curve
plt.title('full sequence E-value ROC curve')
plt.plot(fprA2, tprA2, color='green', label='full sequence E-value')
plt.legend(loc='lower right')
plt.xlim([-0.01, 1])
plt.ylim([0, 1.03])
```

```
plt.ylabel('True positive rate')
plt.xlabel('False positive rate')
plt.show()
aucA2
```



[66]: 0.9999990440426826

```
[35]: X_trainB2 = X_testB
Y_trainB2 = Y_testB
X_testB2 = X_trainB
Y_testB2 = Y_trainB
X_trainB2_2D = X_testB_2D
X_testB2_2D = X_trainB_2D
```

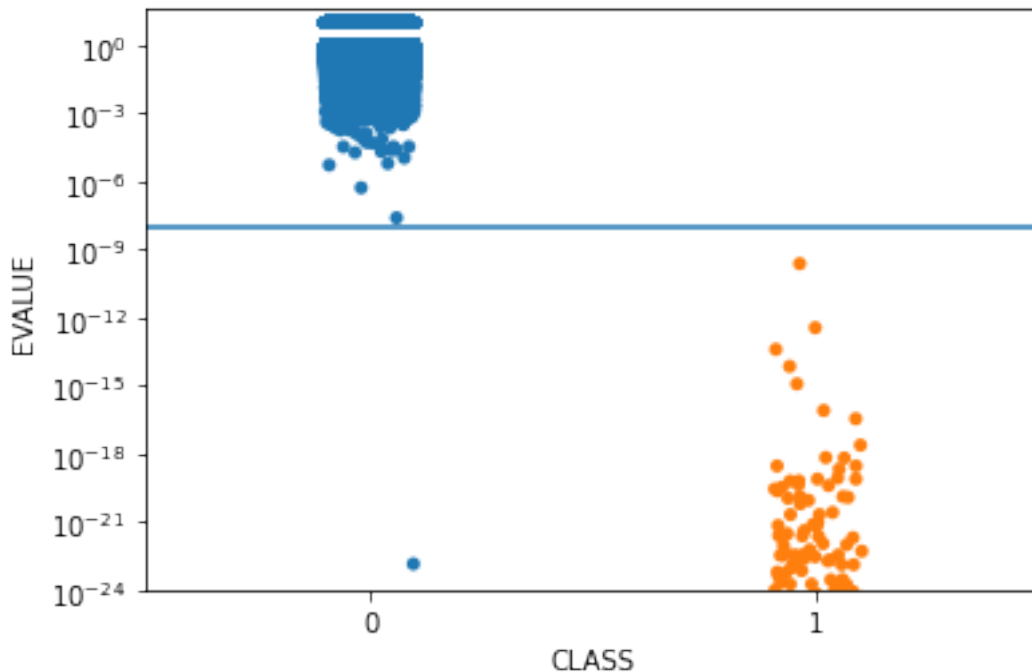
```
[36]: #computation of MCC for best domain evalule for each treshold (ranging from
      ↪ 1e-20 to 1)
mcc_list = []
for i in range(-20, 1, 1):
    Y_predB2 = X_trainB2.apply(lambda x: 1 if x<10**(i) else 0)
    mcc_i = matthews_corrcoef(Y_trainB2, Y_predB2)
    mcc_list.append((i, mcc_i))
    Y_predB2 = Y_predB2.iloc[0:0]
print(mcc_list)
```

```
[(-20, 0.914425489031218), (-19, 0.9581664291145868), (-18, 0.9612147814910796),
(-17, 0.967282730981766), (-16, 0.976314105338485), (-15, 0.9793060974290068),
(-14, 0.9793060974290068), (-13, 0.988227849695091), (-12, 0.9911839636027662),
(-11, 0.9911839636027662), (-10, 0.9941313083026798), (-9, 0.9941313083026798),
(-8, 0.9912012743277198), (-7, 0.9912012743277198), (-6, 0.9941993109664184),
(-5, 0.9884983283299139), (-4, 0.9562355491976475), (-3, 0.8016481568368459),
(-2, 0.4025392054252638), (-1, 0.12364184341982822), (0, 0.039263210351458626)]
```

```
[38]: #best th= 10e-9 best domain evalve
#Th applied on test
Y_predB2 = X_testB2.apply(lambda x: 1 if x<10e-9 else 0)
mismatch = Y_predB2.loc[Y_testB2 != Y_predB2]
false_pos = mismatch.loc[mismatch == 1].index
false_neg = mismatch.loc[mismatch == 0].index
print('false positives are:', false_pos)
print('false negatives are:', false_neg)
striplot = sns.striplot(y = X_testB2, x = Y_testB2)
striplot.set(yscale='log', ylim=(10e-25,40))
striplot.axhline(10e-9)
```

```
false positives are: Index(['G3LH89'], dtype='object')
false negatives are: Index([], dtype='object')
```

```
[38]: <matplotlib.lines.Line2D at 0x7fcc965af0d0>
```



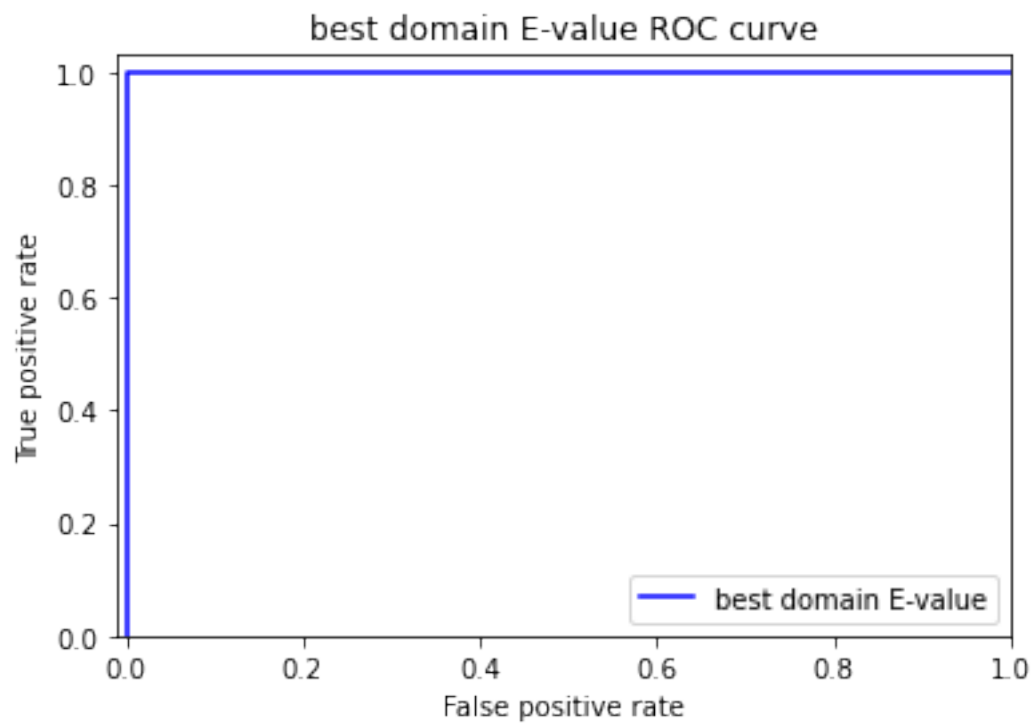
```
[39]: sklearn.metrics.accuracy_score(Y_testB2, Y_predB2)
```

```
[39]: 0.9999964428002277
```

```
[40]: #confusion matrix best domain  
from sklearn.metrics import confusion_matrix  
confusion_matrix(Y_testB2, Y_predB2)
```

```
[40]: array([[280944,    1],  
        [    0,   175]])
```

```
[41]: #computation of ROC curve and ACU for best domain e-value  
import sklearn  
from sklearn import metrics  
import matplotlib.pyplot as plt  
Y_score = [- values for values in X_testB2]  
fprB2, tprB2, thB2 = sklearn.metrics.roc_curve(Y_testB2, Y_score)  
aucB2 = metrics.auc(fprB2, tprB2)  
#plot ROC curve  
plt.title('best domain E-value ROC curve')  
plt.plot(fprB2, tprB2, color='blue', label='best domain E-value')  
plt.legend(loc='lower right')  
plt.xlim([-0.01, 1])  
plt.ylim([0, 1.03])  
plt.ylabel('True positive rate')  
plt.xlabel('False positive rate')  
plt.show()  
aucB2
```



[41]: 0.9999987796289563

[ ]: