

Analysis of Model Evaluation Metrics

David Wang

March 30 2023

1 Precision and Recall

Precision and recall are common metrics used to evaluate the performance of a binary classification model.

We are given a dataset and a model which predicts, for each sample in the dataset, True or False. We define this as the positive class (True) and the negative class (False).

We also define:

- The true positives (TP) as the number of "True" predictions that are actually correct.
- The false positives (FP) as the number of "True" predictions that are not correct.
- The true negatives (TN) as the number of "False" predictions that are correct.
- The false negatives (FN) as the number of "False" predictions that are incorrect.

Now we define

- Precision, as the fraction of "True" predictions made by the model that should actually be "True"

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

- Recall, as the fraction of "True" examples which was correctly identified as "True" by the model

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

However, if we have an unbalanced dataset, then these metrics may not be the most informative, as will be demonstrated below.

2 Matthew's Correlation Coefficient

From some preliminary research, another metric which could be utilized is the Matthews's correlation coefficient (MCC), which is defined as

$$MCC = \frac{TN \cdot TP - FN \cdot FP}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

The score produced by MCC provides a good evaluation on the performance of the model, even with unbalanced priors.

3 Tests

In the case where we care about the performance of the model on both the "True" and "False" labels, it is more intuitive to think of the two labels as "class A" and "class B". If we use precision and recall, the metrics do not take into account the number of elements in each class, and can produce very skewed results in the case of unbalanced priors. A Python script (A.1) is written to test and illustrate a more intuitive understanding of these metrics. For each test, we construct a confusion matrix of the form

$$\begin{bmatrix} TP & FP \\ FN & TN \end{bmatrix}$$

and compare the efficacy of accuracy, precision, recall, and MCC on evaluating the performance of the model. We also swap the labels (ie change the negative class into the positive class) to show the potential problems with evaluation using precision and recall.

3.1 Test 1

```
>>>> test 1 >>>>>
[[999  1]
 [  9  1]]
1000 in class A, 999 correctly predicted,
10 in class B, 1 correctly predicted
accuracy: 0.990, precision: 0.991, recall: 0.999
mcc: 0.220

-----swapping labels-----
10 in class A, 1 correctly predicted,
1000 in class B, 999 correctly predicted
accuracy: 0.990, precision: 0.500, recall: 0.100
mcc: 0.220
```

Initially, we see that accuracy, precision, and recall are all very high. However, when we switch the labels between class A and class B, we notice that the precision and recall goes down significantly. So even though the model had a

9/10 error rate in predicting class B, because the number of samples for class B was so low, this did not reflect with precision and recall. We see, however, that MCC is invariant to the switching of labels and accurately reflects the bad performance of our model in predictions of class B.

3.2 Test 2

```
>>>> test 2 >>>>>
[[999  1]
 [ 1  9]]
1000 in class A, 999 correctly predicted,
10 in class B, 9 correctly predicted
accuracy: 0.998, precision: 0.999, recall: 0.999
mcc: 0.899

-----swapping labels-----
10 in class A, 9 correctly predicted,
1000 in class B, 999 correctly predicted
accuracy: 0.998, precision: 0.900, recall: 0.900
mcc: 0.899
```

Here, per class accuracy is high for A, but only 90% for B. However, the number of samples in B is again very small, and thus, neither precision, accuracy, nor recall reflects the bad performance seen for that class. When the labels are swapped, we now see the drop in precision and recall. Meanwhile, MCC accounts for this without the need of label swapping and re-evaluating.

3.3 Test 3

```
>>>> test 3 >>>>>
[[999  1]
 [998  2]]
1000 in class A, 999 correctly predicted,
1000 in class B, 2 correctly predicted
accuracy: 0.500, precision: 0.500, recall: 0.999
mcc: 0.013

-----swapping labels-----
1000 in class A, 2 correctly predicted,
1000 in class B, 999 correctly predicted
accuracy: 0.500, precision: 0.667, recall: 0.002
mcc: 0.013
```

Here, the dataset contains equal distributions of samples from class A and class B. However, model performance on class B is extremely poor. We note that accuracy is affected since the number of samples in each class is similar. However, the extremely bad performance on class B is only obvious when we

evaluate the recall after switching labels. Meanwhile, MCC is able to accurately reflect the poor performance of the model.

3.4 Test 4

```
>>>> test 4 >>>>>
[[ 1 999]
 [ 9  1]]
1000 in class A, 1 correctly predicted,
10 in class B, 1 correctly predicted
accuracy: 0.002, precision: 0.100, recall: 0.001
mcc: -0.899

-----swapping labels-----
10 in class A, 1 correctly predicted,
1000 in class B, 1 correctly predicted
accuracy: 0.002, precision: 0.001, recall: 0.100
mcc: -0.899
```

Here, we see that the model has bad performance on both classes, and MCC becomes negative in this case.

3.5 Test 5

```
>>>> test 5 >>>>>
[[999  1]
 [ 2 998]]
1000 in class A, 999 correctly predicted,
1000 in class B, 998 correctly predicted
accuracy: 0.999, precision: 0.998, recall: 0.999
mcc: 0.997

-----swapping labels-----
1000 in class A, 998 correctly predicted,
1000 in class B, 999 correctly predicted
accuracy: 0.999, precision: 0.999, recall: 0.998
mcc: 0.997
```

Here, we see that the model has good performance on both classes, and there is an equal number of samples for each class. MCC is high with no surprise here, and precision and recall also correctly evaluate the good-performance of the model.

4 Relation to True and False Positive Rates

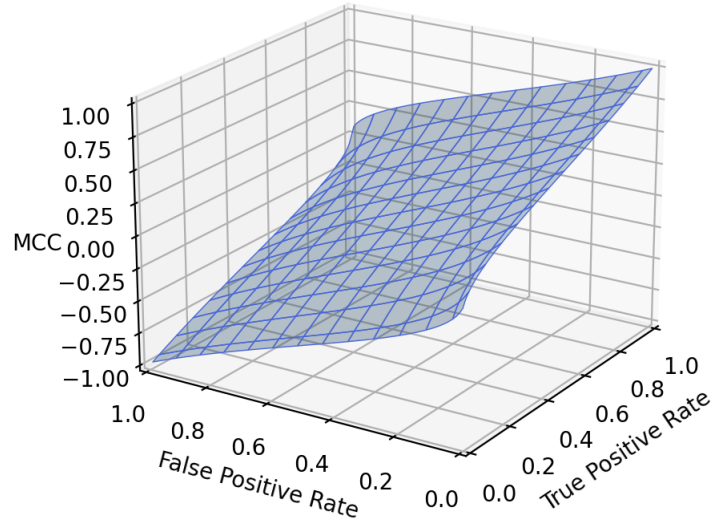
It is interesting to see the relationship between the Matthew's Correlation Coefficient to the true and false positive rates. The true positive rate is the recall

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

and the false positive rate is $1 - \text{recall}$ after the labels have been switched

$$\text{FPR} = \frac{FP}{FP + TN}$$

Writing a script in python (A.2) to visualize the relation, the following plot is obtained



This confirms that MCC can be used as a good metric for evaluating our model, as at a max value of $\text{MCC} = 1.0$, the model performs perfectly with a TPR of 1.0 and FPR of 0.0.

5 Conclusion

From this initial exploration, we determine that the Matthew's Correlation Coefficient provides a more robust method to evaluating the performance of a binary classification model, which is insensitive to the prior distributions. Interesting further explorations could investigate a potential generalization of this metric to multi-class classification.

A Appendix

The following Python scripts were written and utilized to explore the effectiveness of MCC. The codebase can also be found at

https://github.com/davidw0311/eval_metrics

A.1 Code for Test Cases

```
1  import numpy as np
2
3  def get_metrics(CM):
4      tp = CM[0,0]
5      tn = CM[1,1]
6      fp = CM[1,0]
7      fn = CM[0,1]
8
9      # precision, recall, and accuracy
10     precision = tp/(tp + fp)
11     recall = tp/(tp+fn)
12     acc = (tp+tn)/(tp+tn+fp+fn)
13     print(f'accuracy: {acc:.3f}, precision: {precision:.3f}, recall: {recall:.3f}')
14
15     # matthew's correlation coefficient
16     mcc = (tn*tp - fn*fp)/np.sqrt((tp+fp)*(tp+fn)*(tn+fp)*(tn+fn))
17     print(f'mcc: {mcc:.3f}')
18
19 def evaluate(CM):
20     print(f'{CM[0,0]+CM[0,1]} in class A, {CM[0,0]} correctly predicted,')
21     print(f'{CM[1,0]+CM[1,1]} in class B, {CM[1,1]} correctly predicted')
22     get_metrics(CM)
23
24
25 def test(CM, num):
26     print(f'>>>>> test {num} >>>>>')
27     print(CM)
28     evaluate(CM)
29
30     print('\n-----swapping labels-----')
31     # switch labels by swapping rows, then swapping columns
32     CM = CM[[1,0]][:,[1,0]]
33     evaluate(CM)
34     print('\n\n')
35
36
37 # create test cases of confusions matrices
38 cms = [
39     np.array([[999, 1],
40              [9, 1]]),
41     np.array([[999, 1],
42              [1, 9]]),
43     np.array([[999, 1],
44              [998, 2]]),
45     np.array([[999, 1],
46              [2, 998]])
47 ]
48
49 for i, cm in enumerate(cms):
50     test(cm, i+1)
```

A.2 Code for Plots of TPR and FPR to MCC

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  def get_metrics(CM):
5      tp = CM[0,0]
6      tn = CM[1,1]
7      fp = CM[1,0]
8      fn = CM[0,1]
9
10     tpr = tp/(tp+fn)
11     fpr = fp/(fp + tn)
12     mcc = (tn*tp - fn*fp)/np.sqrt((tp+fp)*(tp+fn)*(tn+fp)*(tn+fn))
13
14     return tpr, fpr, mcc
15
16
17  pAs = []
18  grid = 100
19  tprs = np.zeros((grid,grid))
20  fprs = np.zeros((grid,grid))
21  mcs = np.zeros((grid,grid))
22
23  for i in range(grid):
24      for j in range(grid):
25          CM = np.array([[grid-i, i],[grid-j,j]])
26          tpr, fpr, mcc = get_metrics(CM)
27          tprs[i, :] = tpr
28          fprs[:, j] = fpr
29          mcs[i,j] = mcc
30
31
32  ax = plt.figure().add_subplot(projection='3d')
33
34  ax.plot_surface(tprs, fprs, mcs, edgecolor='royalblue', lw=0.5, rstride=8, cstride=8,
35                 alpha=0.3)
36  ax.set(xlim=(0, 1), ylim=(0, 1), zlim=(-1, 1),
37         xlabel='True Positive Rate', ylabel='False Positive Rate', zlabel='MCC')
38
39  plt.show()
```