

Ans 3:

(a)

Relationship between ROC and PR curve is:

We evaluate the performance of machine learning algorithm of ROC and PR curve on same dataset. There are certain theorems which explain the relationship between them.

(Refer from :<https://www.biostat.wisc.edu/~page/rocpr.pdf>)(these 3 the)

-Theorem 3.1:

For a given dataset of positive and negative examples, there exists a one-to-one correspondence between a curve in ROC space and a curve in PR space, such that the curves contain exactly the same confusion matrices, if Recall not equal 0.

*Here one to one correspondence values should be uniquely mapped with one another.

*Suppose Recall is equal to 0 then we can't be able to find out False positive, also we can't be able to find out True negative.

*If recall is equal to 0 then True positive is 0, then if we see the formula of Recall then numerator is 0 so we can't be able to find out False Negative. So if recall not equal to 0 then we can never recover False Positive and hence we can never find out True negative also.

Formulas: Definitions of metrics

$$\text{Recall} = TP / TP + FN$$

$$\text{Precision} = TP / TP + FP$$

$$\text{True Positive Rate} = TP / TP + FN$$

$$\text{False Positive Rate} = FP / FP + TN$$

*hence this one to one mapping between confusion matrices and points in PR space this acknowledge we can translate curve in ROC space to PR space and vice versa.

-Theorem 3.2:

*For a fixed number of positive and negative examples, one curve dominates a second curve in ROC space if and only if the first dominates the second in Precision-Recall space. This explanation is further in 2nd part of this question.

* if one curve dominates in other curve in ROC same thing will be happened in PR curve and vice versa explained in this theorem.

-Corollary:

Let assume you have a set of given points in PR space ,then there is a PR curve plot that dominates the other valid PR curves that could be constructed with these points.(This part will be further explained in part 3 of this question)

*It is observed that ROC would present over optimistic view of algorithm but PR doesn't ,when there is large skewed in class distribution of data points.

* As we are talking about curves differentiation of PR and ROC in between different algorithms then we observe that PR curve can expose difference in algorithms that are not apparent in ROC space.

(b) Prove that a curve dominates in ROC if and only if it dominates in PR curve.

As explained in the paper given, this can be proved by contradiction:

*Claim 1: If a curve dominates in ROC space then it dominates in PR space.

Let suppose there are 2 curves I and II and assume curve I in ROC space dominates while curve I will never dominate when we translate this curve in PR space.

Let assume point A on curve II and point B on curve I with same Recall have lower Precision.

$\text{Precision}(A) > \text{Precision}(B)$ yet $\text{Recall}(A) = \text{Recall}(B)$, Recall is identical to TPR then $\text{TPR}(A) = \text{TPR}(B)$. Since we know curve I dominates curve II in ROC space then $\text{FPR}(A) \geq \text{FPR}(B)$.

(these formulas taken as it is)

As we know that total positives and total negatives are fixed and since $\text{TPR}(A) = \text{TPR}(B)$:
 $\text{TPR}(A) = \text{TPA} / \text{Total Positives}$, $\text{TPR}(B) = \text{TPB} / \text{Total Positive}$, we now have $\text{TPA} = \text{TPB}$ and thus denote both as TP. Remember that $\text{FPR}(A) \geq \text{FPR}(B)$ and

$$\text{FPR}(A) = \text{FPA} / \text{Total Negatives}, \text{FPR}(B) = \text{FPB} / \text{Total Negatives}$$

This implies that $\text{FPA} \geq \text{FPB}$ because, $\text{PRECISION}(A) = \text{TP} / \text{FPA} + \text{TP}$,

$$\text{PRECISION}(B) = \text{TP} / \text{FPB} + \text{TP},$$

we now have that $\text{PRECISION}(A) \leq \text{PRECISION}(B)$. But this contradicts our original assumption that $\text{PRECISION}(A) > \text{PRECISION}(B)$. Hence proved.

*Claim 2:

If a curve dominates in PR space then it dominates in ROC space.

By contradiction let us assume we have 2 curves like in claim 1, curve I and curve II. And assume curve I dominates curve II in PR space but curve I no longer dominates once it is translated into ROC space. As we saw curve I no longer dominates in ROC space then there should exist a point A on Curve II such that point B in curve I both have the same TPR and also;

$FPR(A) < FPR(B)$. Also Recall and TPR are the same then $Recall(A) = Recall(B)$. Also we saw curve I dominates in PR space so $Precision(A) \leq Precision(B)$. As $FP(A) \geq FP(B)$ then $FPR(A) \geq FPR(B)$ that contradicts our assumption that we assumed initially.

$RECALL(A) = \frac{TPA}{\text{Total Positives}}$, $RECALL(B) = \frac{TPB}{\text{Total Positives}}$

As we know $PRECISION(A) \leq PRECISION(B)$ and $PRECISION(A) = \frac{TP}{TP + FPA}$,
 $PRECISION(B) = \frac{TP}{TP + FPB}$

we find that $FPA \geq FPB$.

Now we have $FPR(A) = \frac{FPA}{\text{Total Negatives}}$, $FPR(B) = \frac{FPB}{\text{Total Negatives}}$

This concludes that if a curve dominates in PR space over curve II (consider other curves) then the same thing will happen if translated to ROC space of curves.

(c) It is incorrect to interpolate between points in PR space. When and why does this happen? How will you tackle this problem?

(Refer from : <https://www.biostat.wisc.edu/~page/rocpr.pdf>) (these 2 points taken as it is)

* To achieve the same we need to do linear interpolation.

Convex hull of points in ROC space represents the curve that dominates all the valid curves that can be constructed with the same points. Such a curve is called an achievable PR curve.

* Performing linear interpolation leads to an overly optimistic estimate of the performance.

But doing linear interpolation to achieve the same is incorrect. Because as recall varies, precision does not vary linearly with recall as it has FP instead of FN in the denominator.

* Following are the procedures :

= Problem of getting the achievable PR curve by converting the points of PR space to ROC space by the help of convex hull. First make the convex hull for ROC curve then

translate it to convex hull of PR space for getting the achievable PR curve. This will result in the same as required curve but near to true PR curve.

=Modified interpolation to get intermediate points A and B which are apart used the following values; $FP(A)$, $FP(B)$, $TP(A)$, $TP(B)$. Then new points to get values are created by adding $TP(A)+x$, where x values between 1 to $TP(B)-TP(A)$. And then we get FB by increasing the false positive linearly for each new point by local skew where local skew is formulated as the ratio of $-FP(B)-FP(A)$ and $TP(B)-TP(A)$, which are corresponding to each point of the data. Then intermediate PR points will be created as per given in paper.