

Analysis Assignment-3

Q1 Static Score

As we could see on different values of r that how much to send to high list. Output was behaving in different ways.

It was seen that best outcome came when r is near 20 or 25. As the mean of length of posting list of all terms was near about 20. Considering some False positive we could also run on $r=25$.

As can be seen in given query.

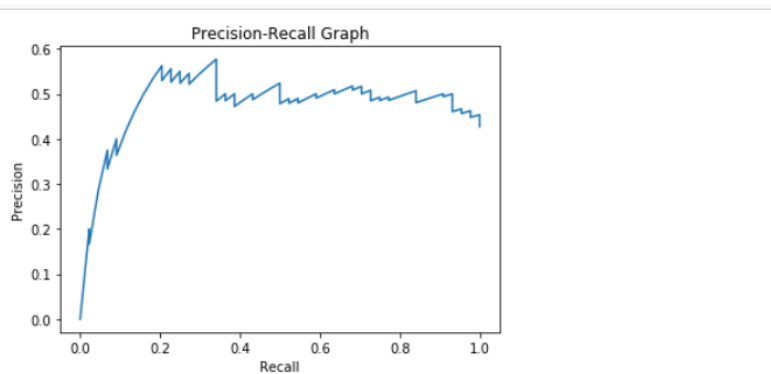
```
querylist=[]
lemmatizer=WordNetLemmatizer()
for word in result:
    if word.isdecimal()==True:
        word=num2words(word)
        querylist.append(lemmatizer.lemmatize(word))
print(querylist)
```

```
I am looking for publically accessible sources of data depicting brain and neuron functions. Especially interesting would be volumetric data from brain scans, electromicographs, and so forth. Please email me and let me know if you know of such.
Number of documents to retrieve
10
['looking', 'publically', 'accessible', 'source', 'data', 'depicting', 'brain', 'neuron', 'function', 'especially', 'interesting', 'would', 'volumetric', 'data', 'brain', 'scan', 'electromicographs', 'forth', 'please', 'email', 'let', 'know', 'know']
```

```
rank(k)
```

```
1008 : 1.0738742724290513
13011 : 1.0669519141006238
13277 : 0.7920841363967563
13012 : 0.5571023136207442
13075 : 0.4853533073379789
651 : 0.41799917574875545
14591 : 0.3754413187145786
1009 : 0.3734288442329891
4533 : 0.3689030243559645
1692 : 0.36439134129337863
```

Here when kept $r=20$.Mean length We are getting better result than $r=15$. Hence this is analysed to be better



As we can infer from graph recall is always increasing with time or its non decreasing in nature and finally reaches 1.

Precision is initially increases and then jitters can be seen .

Hence the relation between precision and recall is analysed from graph above.

Q3

a) Explain the relationship between ROC curve and PR curve.

Relation between Roc curve and PR curve is such that for a given dataset of positive and negative samples , there exist a one to one correspondence b/w ROC curve & PR curve , if Recall > 0.

As point in a ROC defines unique matrix as TPR and FPR consist of all TP,FP,TN and FN for a fixed dataset .In PR we ignore TN ,this might seem like we may have different confusion matrix .But with given other three enteries for fixed data TN is also uniquely identified.This is true only when TPR is not 0.If recall is 0 we cant measure FN.Hence no uniqueness left.

Hence Confusion matrix is in one to one relation with ROC and PR curve.This implies ROC and FP are in one to one relation with each other.

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

We will prove the above in two claims stated below:

Consider two curves curve 1 and curve 2 with point B on curve 1 and point A on curve 2.

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

Proof by Contradiction

Let us assume Curve 1 dominate in ROC Space and do not dominate in PR space.

Since curve 1 does not dominate in PR curve then there will be one point A on curve 2 whose precision is greater than one point B on curve 1 point on same Recall.

i.e $\text{PRECISION}(A) > \text{PRECISION}(B)$ at $\text{TPR}(A) = \text{TPR}(B)$.----eq1

$\text{FPR}(B) = \text{False Positive}(B) / \text{Total Negatives}$

$\text{FPR}(A) = \text{False Positive}(A) / \text{Total Negatives}$

Now curve 1 dominates curve 2 in ROC therefore $\text{FP}(A) \geq \text{FP}(B)$ because Total Negatives are constant.

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

Now here we can see that $\text{PRECISION}(A) \leq \text{PRECISION}(B)$. but this contradicts our assumption(eq1) Therefore Claim 1 is true.

Claim 2 : If a curve dominates in PR space then it dominates in ROC space.

Let us consider that curve 1 dominates in PR space but does not dominate in ROC space.

Proof by Contradiction

AS curve 1 dominates in PR space we know that $\text{PRECISION}(A) \leq \text{PRECISION}(B) \rightarrow (I)$, but

$\text{RECALL}(A) = \text{RECALL}(B)$.

$\text{PRECISION}(A) = \text{TP} / (\text{TP} + \text{FP}(A))$

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

By (I) we can say that $\text{FP}(A) \geq \text{FP}(B)$

Now

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

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

From above equation we can say that $\text{FPR}(A) \geq \text{FPR}(B)$ and this contradicts our original assumption that $\text{FPR}(A) < \text{FPR}(B)$. Hence it also dominates in roc.

Hence proved roc dominates if and only if pr dominates.

Q3:]

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

It is incorrect to interpolate between points in precision recall space because as FP is considered in Precision and FN is considered in Recall and hence in different levels recall varies but it Precision does not necessarily change. In order to tackle this problem, it is better to translate the points to ROC curve(convex hull in ROC space) and again from ROC curve to PR space. Now, this curve PR space excludes exactly those points which are beneath the convex hull in ROC space.

To interpolate between two points A and B, we must interpolate between their counts $\text{TP}(A)$, $\text{TP}(B)$, $\text{FP}(A)$ and $\text{FP}(B)$.

For this we find local skew defined as:

$\text{FP}(B) - \text{FP}(A) / \text{TP}(A) - \text{TP}(B)$

Now, for all x where $1 \leq x \leq \text{TP}(A) - \text{TP}(B)$, we create new points $\text{TP}(A) + x$ and calculate corresponding FP by increasing it linearly by local skew. The intermediate points are:
 $(\text{TP} / \text{Total positives}, (\text{TP}(A) + x) / ((\text{TP}(A) + x) \text{FP}(A) + \text{local skew}))$

Hence the problem can be solved in this way by newly constructed PR space curve suitable for interpolation.

