

```
In [4]: dhtml("Avanish Singh")  
dhtml("191550022")
```

Avanish Singh

191550022

```
In [5]: q = "new new times"  
d1 = "new york times"  
d2 = "new york post"  
d3 = "los angeles times"
```

```
In [6]: def jaccard(Query, Statement):  
    words_doc1 = set(Query.lower().split())  
    words_doc2 = set(Statement.lower().split())  
    intersection = words_doc1.intersection(words_doc2)  
    union = words_doc1.union(words_doc2)  
    return float(len(intersection)) / len(union)
```

```
In [7]: jaccard(q, d1)
```

```
Out[7]: 0.6666666666666666
```

```
In [8]: jaccard(q, d2)
```

```
Out[8]: 0.25
```

```
In [9]: jaccard(q, d3)
```

```
Out[9]: 0.25
```

```
In [11]: import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
```

```
In [12]: def cos_sim(Query, Statement):
    X_list = word_tokenize(Query)
    Y_list = word_tokenize(Statement)
    sw = stopwords.words('english')
    l1=[];l2=[]
    X_set = {w for w in X_list if not w in sw}
    Y_set = {w for w in Y_list if not w in sw}
    rvector = X_set.union(Y_set)
    for w in rvector:
        if w in X_set: l1.append(1)
        else: l1.append(0)
        if w in Y_set: l2.append(1)
        else: l2.append(0)
    c = 0
    for i in range(len(rvector)):
        c+= l1[i]*l2[i]
        cosine = c / float((sum(l1)*sum(l2))**0.5)
    return cosine
```

```
In [13]: cos_sim(q, d1)
```

```
Out[13]: 0.8164965809277261
```

```
In [14]: cos_sim(q, d2)
```

```
Out[14]: 0.4082482904638631
```

```
In [15]: cos_sim(q, d3)
```

```
Out[15]: 0.4082482904638631
```

```
In [ ]:
```

