

Project for the Course: Meaningful Predictive Modeling.

Description of the Dataset:

The dataset that will be analysed for this project has been downloaded from the Kaggle website. The dataset - 'Flipcart Products Review Dataset' Description: The dataset contains information such as products name, price, review, rating and summary for 104 types of different products on flipcart.com. It has 194276 rows and 5 columns. It was chosen for analysis as it has different customer reviews and can be easily used for sentiment analysis.

To download the data and for more description on the datasets go to: <https://www.kaggle.com/datasets/mansithummar67/flipkart-product-review-dataset?resource=download> (<https://www.kaggle.com/datasets/mansithummar67/flipkart-product-review-dataset?resource=download>)

Objectives:

1. Whether the product reviews on a product can be used to determine overall sentiment towards the product.
2. Whether the product reviews can be used for pricing decisions. That is, if a product has higher ratings and positive reviews, the price may be increased, if it has lower ratings and negative reviews, the price may be decreased.

In [1]:

```
1
2 import gzip
3 import string
4 import random
5 from nltk.stem.porter import PorterStemmer
6
7 path = r"D:\TRL\OneDrive\CourseraPython\Final_Course3\flipkart_product.csv"
```

The following code can be used to get the data from the csv file BUT: see the note below.

Note: The following code was first used to read the data from the csv file, but it does not work the way that it was intended to because the csv file as well as the Product descriptions in the first column are comma separated. Because of this the product description spills over to the next column instead of the next column entry. Also using any other delimiter does not help with this problem. So I have researched pandas and used them.

```
f = open(path,'rt', encoding = "ISO-8859-1")

import csv
reader = csv.reader(f, delimiter = ',')

header = f.readline()
header = header.strip().split(',')
header = header

d = next(reader) # there is some problem with the csv file, Prices in all d # entries have garbage before the price value, which I have to try and remove

dataset = []
for line in f:
    fields = line.strip().split("\t")
    d = dict(zip(header,fields)) # keys in the header are mapped for field in ['Price', 'Rate']: # to values in each line
    d[field] = int(d[field])
    dataset.append(d)

print("len(d) =",len(dataset))
d

dataset[10] # testing

dataset[189871]
```

The following also doesn't work:

```
data_new = []
for lines in range(len(dataset)):
    data_new.append(lines)

type(dataset)

type(data_new)
```

In [2]:

```
1 ### using pandas and dataframe to read the csv file as it cannot be
2 # read as dictionary object using a for loop
3
4 import pandas as pd
5 df = pd.read_csv("D:\TRL\OneDrive\CourseraPython\Final_Course3\flipkart_product.csv")
```

In [3]:

```
1 import chardet
2 with open(path, 'rb') as rawdata:
3     result = chardet.detect(rawdata.read(189872))
```

In [4]:

```
1 result
```

Out[4]:

```
{'encoding': 'Windows-1254',
 'confidence': 0.4569690721335204,
 'language': 'Turkish'}
```

In []:

```
1 dataset = pd.read_csv(r"D:\TRL\OneDrive\CourseraPython\Final_Course3\flipkart_product.csv", encoding ="ISO-8859-1",dtype = {'ProductN
```

In [6]:

```
1 dataset
```

Out[6]:

	ProductName	Price	Rate	Review	Summary
0	Candes 12 L Room/Personal Air Cooler?ÿ?ÿ(White...	3999	5	Super!	Great cooler.. excellent air flow and for this...
1	Candes 12 L Room/Personal Air Cooler?ÿ?ÿ(White...	3999	5	Awesome	Best budget 2 fit cooler. Nice cooling
2	Candes 12 L Room/Personal Air Cooler?ÿ?ÿ(White...	3999	3	Fair	The quality is good but the power of air is de...
3	Candes 12 L Room/Personal Air Cooler?ÿ?ÿ(White...	3999	1	Useless product	Very bad product it's a only a fan
4	Candes 12 L Room/Personal Air Cooler?ÿ?ÿ(White...	3999	3	Fair	Ok ok product
...
189868	NIVEA Soft Light Moisturizer for Face, Hand & ...	142	5	Terrific	Good
189869	NIVEA Soft Light Moisturizer for Face, Hand & ...	142	5	Terrific	Thanks
189870	NIVEA Soft Light Moisturizer for Face, Hand & ...	142	5	Terrific	Good
189871	NIVEA Soft Light Moisturizer for Face, Hand & ...	142	5	Just wow!	Super
189872	NIVEA Soft Light Moisturizer for Face, Hand & ...	142	4	Worth the money	Good

189873 rows × 5 columns

In [7]:

```
1 dataset.iloc[1,:]
```

Out[7]:

ProductName Candes 12 L Room/Personal Air Cooler?ÿ?ÿ(White...
Price 3999
Rate 5
Review Awesome
Summary Best budget 2 fit cooler. Nice cooling
Name: 1, dtype: object

In [8]:

```
1 dataset.iloc[1,2]
```

Out[8]:

'5'

In [9]:

```
1 from collections import defaultdict  
2 import numpy as np
```

In [10]:

```
1 # for the sake of analysis, let us randomly generate a column of number of  
2 # helpful reviews, with the number of helpful reviews between 1 to 100.  
3  
4 dataset['Helpful_Reviews'] = np.random.randint(1,100,dataset.shape[0])
```

In [11]:

```
1 dataset
```

Out[11]:

	ProductName	Price	Rate	Review	Summary	Helpful_Reviews
0	Candes 12 L Room/Personal Air Cooler?ý?ý(White...	3999	5	Super!	Great cooler.. excellent air flow and for this...	61
1	Candes 12 L Room/Personal Air Cooler?ý?ý(White...	3999	5	Awesome	Best budget 2 fit cooler. Nice cooling	49
2	Candes 12 L Room/Personal Air Cooler?ý?ý(White...	3999	3	Fair	The quality is good but the power of air is de...	45
3	Candes 12 L Room/Personal Air Cooler?ý?ý(White...	3999	1	Useless product	Very bad product it's a only a fan	5
4	Candes 12 L Room/Personal Air Cooler?ý?ý(White...	3999	3	Fair	Ok ok product	57
...
189868	NIVEA Soft Light Moisturizer for Face, Hand & ...	142	5	Terrific	Good	36
189869	NIVEA Soft Light Moisturizer for Face, Hand & ...	142	5	Terrific	Thanks	12
189870	NIVEA Soft Light Moisturizer for Face, Hand & ...	142	5	Terrific	Good	51
189871	NIVEA Soft Light Moisturizer for Face, Hand & ...	142	5	Just wow!	Super	22
189872	NIVEA Soft Light Moisturizer for Face, Hand & ...	142	4	Worth the money	Good	64

189873 rows × 6 columns

In [12]:

```
1 # since the last row is NaN, we delete it
2 # dataset.drop(189873) # no need to run it the second time
```

In [13]:

```
1 dataset.iloc[1][1] + dataset.iloc[2][1], dataset.iloc[1][2] + dataset.iloc[2][2], dataset.iloc[1][5] + dataset.iloc[2][5]
2
3 # this shows that the values in the second column and third column are
4 # strings, and on adding, will be concatenated, the values in the sixth
5 # column can be added etc
```

Out[13]:

('39993999', '53', 94)

In [14]:

```
1 int(dataset.iloc[1][1]) + int(dataset.iloc[2][1]) # we have to convert to add
```

Out[14]:

7998

In [17]:

```
1 indices = []
2 for i in range(len(dataset)):
3     if len(str(dataset.iloc[i][2])) > 2 :
4         indices.append(i)
5         print(i)
```

40725
82709
82720

In [18]:

```
1 indices
2 dataset.iloc[indices[0]][2], dataset.iloc[indices[1]][2], dataset.iloc[indices[2]][2]
```

Out[18]:

('Pigeon Favourite Electric Kettle?ý?ý(1.5 L, Silver, Black)',
'Bajaj DX 2 L/W Dry Iron',
'Nova Plus Amaze NI 10 1100 W Dry Iron?ý?ý(Grey & Turquoise)')

In [19]:

```
1 # these three cells were deleted for analyses as no number was present
2 # in the Price and rating fields.
3 dataset.drop([indices[0],indices[1],indices[2]], inplace = True) # inplace = True updates the dataframe
```

In [20]:

```
1 dataset.iloc[indices[2]] # this shows that the rows corresponding to these
2 # indices have been deleted
```

Out[20]:

```
ProductName      Nova Plus Amaze NI 10 1100 W Dry Iron?ý?ý(Grey...
Price                                                    494
Rate                                                    5
Review                                                    Terrific
Summary          Fantastic look, Fast producing heat, Light wei...
Helpful_Reviews                                     38
Name: 82723, dtype: object
```

In [21]:

```
1 len(dataset)
```

Out[21]:

189870

In [22]:

```
1 # Our goal is to build a classifier that estimates the sentiment- the star
2 # rating, based on the occurrence of words in the review.
3 # Note: the column 4 of the dataset has been derived using the Summary.
4
5 # For this, we first count the number of unique words in the reviews
6 wordCount = defaultdict(int)
7 for i in range(len(dataset)):
8     for w in str(dataset.iloc[i][4]).split(): # see code in cells below for how this works
9         wordCount[w] += 1
10
11 print(len(wordCount))
12
```

82191

In [23]:

```
1 type(dataset['Summary'][2]) ,dataset['Summary'][3]
```

Out[23]:

(str, "Very bad product it's a only a fan")

In [24]:

```
1 dataset.iloc[2][4], dataset.iloc[3][4]
```

Out[24]:

('The quality is good but the power of air is decent',
"Very bad product it's a only a fan")

In [25]:

```
1 len(dataset)
```

Out[25]:

189870

In [26]:

```
1 dataset['Summary'][5].split()
```

Out[26]:

```
['The',  
'cooler',  
'is',  
'really',  
'fantastic',  
'and',  
'provides',  
'good',  
'air',  
'flow.',  
'Highly',  
'recommended']
```

In [27]:

```
1 range(len(dataset))
```

Out[27]:

range(0, 189870)

In [28]:

```
1 len(wordCount)
```

Out[28]:

82191

In [29]:

```
1 # Now Let us ignore the capitalization and punctuation
2 wordCount = defaultdict(int)
3 punctuation = set(string.punctuation)
4
5 for i in range(len(dataset)): #dataset['Summary'][i] can be used instead of dataset['Summary'][i]
6     s = ''.join([c for c in str(dataset.iloc[i][4]).lower() if not c in punctuation])
7     for w in s.split():
8         wordCount[w] += 1
9
10 print(len(wordCount))
11 # it did reduce by about 30000
```

49684

In [30]:

```
1 # Now Let us try stemming to see the unique stems of the words in the reviews
2 wordCount = defaultdict(int)
3 punctuation = set(string.punctuation)
4 stemmer = PorterStemmer()
5
6 for i in range(len(dataset)):
7     s = ''.join([c for c in str(dataset.iloc[i][4]).lower() if not c in punctuation])
8     for w in s.split():
9         w = stemmer.stem(w)
10        wordCount[w] += 1
11
12 print(len(wordCount))
13 # it further reduced
```

43318

In [31]:

```
1 # Now Let us extract the most common words
2 count_unique = [(wordCount[w],w) for w in wordCount]
3 count_unique.sort()
4 count_unique.reverse()
5 my_words = [w[1] for w in count_unique[:1000]]
```

In [32]:

```
1 my_words
```

Out[32]:

```
['good',
 'product',
 'is',
 'veri',
 'it',
 'nice',
 'and',
 'the',
 'for',
 'qualiti',
 'thi',
 'i',
 'to',
 'not',
 'in',
 'of',
 'a',
 'but']
```

In [33]:

```
1 my_words_id = dict(zip(my_words,range(len(my_words))))
```

In [34]:

```
1 my_words_id['best'] # example
```

Out[34]:

18

In [35]:

```
1 my_words_set = set(my_words)
```

In [36]:

```
1 # Now these words correspond to feature indices
2 # i.e. each review can be mapped to a 1000 dimensional vector
3
4 def feature(datarow):
5     feat = [0]*len(my_words)
6     s = ''.join([c for c in str(dataset.iloc[datarow][4]).lower() if not c in punctuation])
7     for w in s.split():
8         if w in my_words:
9             feat[my_words_id[w]] += 1
10    feat.append(1) # this is for the offset feature
11    return feat
```

In [37]:

```
1 # Now, for each review, we can calculate rating as follows:
2 # rating = alpha + (sum over word) count(word) * theta(word)
3 # theta(word) is the amount by which rating is affected for the 'word'
4 # in the review.
5
6 # Now to fit this model, we first shuffle the dataset
7 # instead of the function shuffle, we use the fucntion sample for dataframes
8 dataset = dataset.sample(frac=1)
9 dataset
```

Out[37]:

	ProductName	Price	Rate	Review	Summary	Helpful_Reviews
89353	MILTON Thermosteel Flip Lid 500 ml Flask	648	5	Excellent	Its very good	99
50728	Saral Home Jute Door Mat	484	5	Great product	Good	83
185636	Butterfly Rhino Plus Wet Grinder	5299	4	Very Good	Best in the price	72
56843	Men Black Sandal	299	5	Nan	super	75
86391	HARIOM ALL IN ONE Solid Wood Study Table	549	3	Good	Table stand is not fit	18
...
61576	Google Nest Mini (2nd Gen) with Google Assista...	3499	1	Unsatisfactory	Very bad product	89
13273	Hindware Nevio Plus 60 Auto Clean Wall Mounte...	13999	5	Brilliant	Great product, timely delivered. Thanks Flipka...	49
7144	Cosito 144 TC Cotton Double Floral Flat Bedshe...	339	1	Worst experience ever!	Not good because it's quilaty is worst not goo...	99
181132	BAJAJ 3 L Instant Water Geyser (Ivora/ flora, ...	2599	5	Fabulous!	Good product	84
20062	SportSoul Cotton Gym & Athletic Abdomen Suppor...	249	5	Excellent	Nice product and comfortable also	76

189870 rows × 6 columns

In [38]:

```
1 dataset.iloc[0] # see how it has changed now
```

Out[38]:

ProductName MILTON Thermosteel Flip Lid 500 ml Flask
Price 648
Rate 5
Review Excellent
Summary Its very good
Helpful_Reviews 99
Name: 89353, dtype: object

In [39]:

```
1 X = [feature(i) for i in range(len(dataset))] # the feature matrix
2 X[1]
```

Out[39]:

```
[1,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0,
 0.]
```

In [40]:

```
1 dataset['Summary'][1] # this is the 2nd summary in the original dataset
```

Out[40]:

```
'Best budget 2 fit cooler. Nice cooling'
```

In [41]:

```
1 dataset.iloc[1][4] # this is the summary in the shuffled dataset.
2 # Note that the 1s in X[1] correspond to the review below.
```

Out[41]:

```
'Good'
```

In [42]:

```
1 # dimensions of the matrix X:
2 len(X), len(X[0]), len(dataset)
```

Out[42]:

```
(189870, 1001, 189870)
```

In [43]:

```
1
```

In [44]:

```
1 dataset.iloc[98411]
```

Out[44]:

```
ProductName      cello Pack of 18 Opalware Cello Dazzle Lush Fi...
Price              1299
Rate                3
Review              Good
Summary            quality is good but i found one broken bowl. p...
Helpful_Reviews    40
Name: 169875, dtype: object
```

In [45]:

```
1 y = [int(dataset.iloc[i][2]) for i in range(len(dataset))] #extracting the labels for each datapoint
```

In [46]:

```
1 type(y), type(y[1])
```

Out[46]:

```
(list, int)
```

In [47]:

```
1 type(dataset['Price'][1])
```

Out[47]:

```
str
```

In [48]:

```
1 import numpy
2 theta, residuals, rank, s = numpy.linalg.lstsq(X,y, rcond= -1)
```

In [49]:

```
1 theta
```

Out[49]:

```
array([[ 0.31935216,  0.07326229, -0.01697148, ...,  0.04540193,
         0.          ,  4.04208313]])
```

In [50]:

```
1 len(theta)
```

Out[50]:

```
1001
```

In [52]:

```
1 # words that have the most positive or negative sentiment:
2 wordSentiments = list(zip(theta, my_words + ['offset']))
3 wordSentiments.sort()
```

In [54]:

```
1 vec = my_words + ['offset']
2 len(vec), vec[1000]
```

Out[54]:

```
(1001, 'offset')
```

In [55]:

```
1 wordSentiments[1:]
```

Out[55]:

```
[(-158254462064.39267, 'littl'),
 (-119127660731.73515, 'beauti'),
 (-64329410242.87878, 'experi'),
 (-29062485098.33034, 'deliv'),
 (-19866341340.955112, 'receiv'),
 (-14808177446.282583, 'materi'),
 (-2367638396.4063454, 'properli'),
 (-1915817829.1457095, 'fridg'),
 (-97958114.31251046, 'cycl'),
 (-82547846.8086258, 'everyth'),
 (-40678272.107992135, 'fantast'),
 (-2067402.7352560484, 'tabl'),
 (-1797726.353305803, 'googl'),
 (-59607.25697928499, 'usag'),
 (-16412.507813690943, 'assembl'),
 (-15819.51320029821, 'sinc'),
 (-8372.694623059942, 'easili'),
 (-6442.748911438002, 'chang').
```

In [56]:

```
1 wordSentiments[-10:]
```

Out[56]:

```
[(248868.8922112999, 'compani'),
 (581033.7833893094, 'condit'),
 (977221.1691484333, 'heavi'),
 (136463205.74634764, 'devic'),
 (601973899.2965547, 'becaus'),
 (1915817829.4525018, 'nois'),
 (2237730774.379421, 'batteri'),
 (76453055425.66623, 'servic'),
 (82040375053.65384, 'ani'),
 (116879229366.5479, 'realli)]
```

In [57]:

```
1 from sklearn import linear_model
```


In [58]:

```
1 help(linear_model.Ridge)
```

Help on class Ridge in module sklearn.linear_model._ridge:

```
class Ridge(sklearn.base.MultiOutputMixin, sklearn.base.RegressorMixin, _BaseRidge)
| Ridge(alpha=1.0, *, fit_intercept=True, normalize='deprecated', copy_X=True, max_iter=None, tol=0.001, solver='auto',
| positive=False, random_state=None)
|
| Linear least squares with l2 regularization.
|
| Minimizes the objective function::
|
|  $\|y - Xw\|_2^2 + \alpha \|w\|_2^2$ 
|
| This model solves a regression model where the loss function is
| the linear least squares function and regularization is given by
| the l2-norm. Also known as Ridge Regression or Tikhonov regularization.
| This estimator has built-in support for multi-variate regression
| (i.e., when y is a 2d-array of shape (n_samples, n_targets)).
|
| Read more in the :ref:`User Guide <ridge_regression>`.
```

In [59]:

```
1 # Let us try ridge regression now with regularization strength 3
2
3 model = linear_model.Ridge(3.0, fit_intercept = False)
4 model.fit(X,y)
```

Out[59]:

Ridge(alpha=3.0, fit_intercept=False)

In [60]:

```
1 theta = model.coef_
```

In [62]:

```
1 wordSentimentsNew = list(zip(theta, my_words + ['offset']))
2 wordSentimentsNew.sort()
```

In [63]:

```
1 wordSentimentsNew[:10]
```

Out[63]:

```
[(-1.8727821610012447, 'poor'),
 (-1.801633747380397, 'wrost'),
 (-1.6433894486160834, 'worst'),
 (-1.6345718876250876, 'bad'),
 (-1.4576261938967328, 'useless'),
 (-1.4266942559614328, 'wast'),
 (-1.3527711098543473, 'third'),
 (-1.335546347037573, 'bed'),
 (-1.0944238213126245, 'broken'),
 (-1.0680260658979142, 'hate')]
```

In [64]:

```
1 # MSE and R2 statistic
2
3 RRprediction = model.predict(X) # prediction using ridge regression
```

In [66]:

```
1 RRprediction, len(RRprediction)
```

Out[66]:

```
(array([4.36117795, 4.36117795, 4.51824222, ..., 0.90841966, 4.43444389,
        4.61213638]),
 189870)
```

In [65]:

```
1 diff = [(a-b)**2 for (a,b) in zip(RRprediction,y)]
```

In [68]:

```
1 MSE = sum(diff)/len(diff)
2 print("MSE = " + str(MSE))
```

MSE = 0.9225290739655049

In [70]:

```
1 FVU = MSE/ numpy.var(y)
2 FVU
```

Out[70]:

0.5418077253507845

In [73]:

```
1 R2 = 1-FVU
2 print("R2 = " + str(R2))
```

R2 = 0.4581922746492155

In [78]:

```
1 # Let us look at this as a classification problem
2 yCategorical = [(rating > 3) for rating in y]
```

In [81]:

```
1 model = linear_model.LogisticRegression(solver = 'lbfgs', max_iter = 500)
2 model.fit(X,yCategorical)
3 # had to increase the number of iterations as the solver did not converge in 100 iterations
```

Out[81]:

LogisticRegression(max_iter=500)

In [84]:

```
1 # Computation of Accuracy of the classifier
2 predictions = model.predict(X)
3 correct = predictions == yCategorical
```

In [86]:

```
1 accuracy = sum(correct)/ len(correct)
2 accuracy
```

Out[86]:

0.8995839258439985

In [87]:

```
1 # True Positive: Prediction is True and Label is True
2 TP = sum([(p and 1) for (p,l) in zip(predictions, yCategorical)]) # 1 and 1 = 1 etc
3 # False Positive: the prediction is True when the Label is False.
4 FP = sum([(p and not 1) for (p,l) in zip(predictions, yCategorical)]) # 1 and not 1 = 0
5 # True negative: prediction is False when actual Label is False.
6 TN = sum([(not p and not 1) for (p,l) in zip(predictions, yCategorical)]) # not 1 and not 1 = 0
7 # False Negative: prediction is False when the actual Label is True
8 FN = sum([(not p and 1) for (p,l) in zip(predictions, yCategorical)]) # not 1 and 1 = 0
9
```

In [88]:

```
1 TPR = TP/(TP+FN)
2 TNR = TN/(TP+FN)
```

In [90]:

```
1 BER = 1-1/2*(TPR+TNR)
2 BER # Balanced error rate
```

Out[90]:

0.42431310162590663

In [92]:

```
1 precision = TP/ (TP+FP)
2 precision
```

Out[92]:

0.9021463232549459

In [93]:

```
1 recall = TP / (TP+FN)
2 recall
```

Out[93]:

0.9775055949524092

In [94]:

```
1 F1 = 2*((precision*recall)/ (precision+recall))
```

In [95]:

```
1 F1
```

Out[95]:

```
0.9383153017910758
```

In [96]:

```
1 # Confidence scores per sample:
2 confidences = model.decision_function(X)
```

In [97]:

```
1 confidences
```

Out[97]:

```
array([ 2.11753817,  2.11753817,  3.27800545, ..., -6.77942346,
        2.42748215,  3.40510561])
```

In [98]:

```
1 len(confidences)
```

Out[98]:

```
189870
```

In [100]:

```
1 confidenceLabels = list(zip(confidences,yCategorical))
```

In [101]:

```
1 confidenceLabels
```

Out[101]:

```
[(2.1175381722499567, True),
 (2.1175381722499567, True),
 (3.278005448773111, True),
 (3.112194695039198, True),
 (-1.9165031678970417, False),
 (7.5066088812710685, True),
 (3.324861824331816, True),
 (4.479214485763893, True),
 (2.1175381722499567, True),
 (-2.955657746575442, False),
 (2.1175381722499567, False),
 (2.847160086211357, True),
 (2.428779031464267, True),
 (2.01084309481325, True),
 (7.441380168165194, True),
 (2.1175381722499567, True),
 (-0.45650721858336285, True),
 (2.524534701165583, True)].
```

In [102]:

```
1 confidenceLabels.sort(reverse = True)
```

In [103]:

```
1 # Labels ranked by confidence
2 labelsRanked = [z[1] for z in confidenceLabels]
```

In [118]:

```
1 def precisionAtK(K,sortedList):
2     return sum(sortedList[:K])/K
```

In [119]:

```
1 def recallAtK(K,sortedList):
2     return sum(sortedList[:K]) / sum(sortedList)
```

In [120]:

```
1 precisionAtK(100,labelsRanked)
```

Out[120]:

```
0.98
```

In [121]:

```
1 precisionAtK(1000,labelsRanked)
```

Out[121]:

0.992

In [122]:

```
1 precisionAtK(10000,labelsRanked)
```

Out[122]:

0.9802

In [123]:

```
1 precisionAtK(100000,labelsRanked)
```

Out[123]:

0.95238

In [124]:

```
1 recallAtK(100,labelsRanked)
```

Out[124]:

0.0006606088386766252

In [125]:

```
1 recallAtK(1000,labelsRanked)
```

Out[125]:

0.006686979264971553

In [126]:

```
1 recallAtK(10000,labelsRanked)
```

Out[126]:

0.06607436568069674

In [127]:

```
1 recallAtK(100000,labelsRanked)
```

Out[127]:

0.6419904548763717

In [128]:

```
1 # Validation pipeline:
2 N = len(X)
3 X_train = X[:N//2]
4 X_valid = X[N//2 : 3*N//4]
5 X_test = X[3*N//4:]
6 y_train = y[:N//2]
7 y_valid = y[N//2:3*N//4]
8 y_test = y[3*N//4:]
```

In [129]:

```
1 len(X), len(X_train), len(X_valid), len(X_test)
```

Out[129]:

(189870, 94935, 47467, 47468)

In [130]:

```
1 # To calculate MSE for a particular model:
2 def MSE(model, X, y):
3     predictions = model.predict(X)
4     differences = [(a-b)**2 for (a,b) in zip(predictions,y)]
5     return sum(differences) / len(differences)
```

In [132]:

```
1 # Train the model for different regularization parameters
2
3 bestModel = None
4 bestMSE = None
5
6 for lamb in [0.01, 0.1, 1, 10, 100]:
7     model = linear_model.Ridge(lamb, fit_intercept = False) # fit a model for each lambda value
8     model.fit(X_train, y_train)
9
10    mseTrain = MSE(model, X_train, y_train)
11    mseValid = MSE(model, X_valid, y_valid)
12
13    print("lambda = " + str(lamb) + ", training/validation error = " +
14          str(mseTrain) + '/' + str(mseValid)) # report the training and validation error
15
16    if not bestModel or mseValid < bestMSE:
17        bestModel = model
18        bestMSE = mseValid
```

lambda = 0.01, training/validation error = 0.9192329263206979/0.9462662313693249

lambda = 0.1, training/validation error = 0.9192425105472739/0.9459396655838058

lambda = 1, training/validation error = 0.9195268042641338/0.9443912918642349

lambda = 10, training/validation error = 0.9208521139900501/0.9414169466619998

lambda = 100, training/validation error = 0.9319774227720088/0.9418901298269216

In [134]:

```
1 # best Model with Least error
2 mseTest = MSE(bestModel, X_test, y_test)
3 print("test error = " + str(mseTest))
4 bestModel
```

test error = 0.9323989084655248

Out[134]:

Ridge(alpha=10, fit_intercept=False)

In [135]:

```
1 bestMSE
```

Out[135]:

0.9414169466619998

In []:

```
1
```