

```
In [ ]: import json
import gzip
import math
from collections import defaultdict
import numpy
from sklearn import linear_model

In [ ]: # This will suppress any warnings, comment out if you'd like to preserve them
import warnings
warnings.filterwarnings("ignore")

In [ ]: # Check formatting of submissions
def assertFloat(x):
    assert type(float(x)) == float

def assertFloatList(items, N):
    assert len(items) == N
    assert [type(float(x)) for x in items] == [float]*N

In [ ]: answers = {}

In [ ]: f = open("spoilers.json.gz", 'r')

In [ ]: dataset = []
for l in f:
    d = eval(l)
    dataset.append(d)

In [ ]: f.close()

In [ ]: # A few utility data structures
reviewsPerUser = defaultdict(list)
reviewsPerItem = defaultdict(list)

for d in dataset:
    u, i = d['user_id'], d['book_id']
    reviewsPerUser[u].append(d)
    reviewsPerItem[i].append(d)

# Sort reviews per user by timestamp
for u in reviewsPerUser:
    reviewsPerUser[u].sort(key=lambda x: x['timestamp'])

# Same for reviews per item
for i in reviewsPerItem:
    reviewsPerItem[i].sort(key=lambda x: x['timestamp'])

In [ ]: # E.g. reviews for this user are sorted from earliest to most recent
[d['timestamp'] for d in reviewsPerUser['b0d7e561ca59e313b728dc30a5b1862e']]

In [ ]: ### 1a

In [ ]:

In [ ]: answers['Q1a'] = MSE(y,ypred)

In [ ]: assertFloat(answers['Q1a'])

In [ ]: ### 1b

In [ ]:

In [ ]: answers['Q1b'] = MSE(y,ypred)

In [ ]: assertFloat(answers['Q1b'])

In [ ]: ### 2

In [ ]: answers['Q2'] = []

for N in [1,2,3]:
    # etc.
    answers['Q2'].append(MSE(y,ypred))

In [ ]: assertFloatList(answers['Q2'], 3)

In [ ]: ### 3a

In [ ]: def feature3(N, u): # For a user u and a window size of N

In [ ]: answers['Q3a'] = [feature3(2,dataset[0]['user_id']), feature3(3,dataset[0]['user_id'])]

In [ ]: assert len(answers['Q3a']) == 2
assert len(answers['Q3a'][0]) == 3
assert len(answers['Q3a'][1]) == 4

In [ ]: ### 3b

In [ ]: answers['Q3b'] = []

for N in [1,2,3]:
    # etc.
    answers['Q3b'].append(mse)

In [ ]: assertFloatList(answers['Q3b'], 3)

In [ ]: ### 4a

In [ ]: globalAverage = [d['rating'] for d in dataset]
globalAverage = sum(globalAverage) / len(globalAverage)

In [ ]: def featureMeanValue(N, u): # For a user u and a window size of N

In [ ]: def featureMissingValue(N, u):

In [ ]: answers['Q4a'] = [featureMeanValue(10, dataset[0]['user_id']), featureMissingValue(10, dataset[0]['user_id'])]

In [ ]: assert len(answers['Q4a']) == 2
assert len(answers['Q4a'][0]) == 11
assert len(answers['Q4a'][1]) == 21

In [ ]: ### 4b

In [ ]: answers['Q4b'] = []

for featFunc in [featureMeanValue, featureMissingValue]:
    # etc.
    answers['Q4b'].append(mse)

In [ ]: assertFloatList(answers["Q4b"], 2)

In [ ]: ### 5

In [ ]: def feature5(sentence):

In [ ]: y = []
X = []

for d in dataset:
    for spoiler,sentence in d['review_sentences']:
        X.append(feature5(sentence))
        y.append(spoiler)

In [ ]:

In [ ]: answers['Q5a'] = X[0]

In [ ]: answers['Q5b'] = [TP, TN, FP, FN, BER]

In [ ]: assert len(answers['Q5a']) == 4
assertFloatList(answers['Q5b'], 5)

In [ ]: ### 6

In [ ]: def feature6(review):

In [ ]: y = []
X = []

for d in dataset:
    sentences = d['review_sentences']
    if len(sentences) < 6: continue
    X.append(feature6(d))
    y.append(sentences[5][0])

#etc.

In [ ]: answers['Q6a'] = X[0]

In [ ]: answers['Q6b'] = BER

In [ ]: assert len(answers['Q6a']) == 9
assertFloat(answers['Q6b'])

In [ ]: ### 7

In [ ]: # 50/25/25% train/valid/test split
Xtrain, Xvalid, Xtest = X[:len(X)//2], X[len(X)//2:(3*len(X))//4], X[(3*len(X))//4:]
ytrain, yvalid, ytest = y[:len(X)//2], y[len(X)//2:(3*len(X))//4], y[(3*len(X))//4:]

In [ ]: for c in [0.01, 0.1, 1, 10, 100]:
    # etc.

In [ ]: answers['Q7'] = bers + [bestC] + [ber]

In [ ]: assertFloatList(answers['Q7'], 7)

In [ ]: ### 8

In [ ]: def Jaccard(s1, s2):
    numer = len(s1.intersection(s2))
    denom = len(s1.union(s2))
    if denom == 0:
        return 0
    return numer / denom

In [ ]: # 75/25% train/test split
dataTrain = dataset[:15000]
dataTest = dataset[15000:]

In [ ]: # A few utilities
itemAverages = defaultdict(list)
ratingMean = []

for d in dataTrain:
    itemAverages[d['book_id']].append(d['rating'])
    ratingMean.append(d['rating'])

for i in itemAverages:
    itemAverages[i] = sum(itemAverages[i]) / len(itemAverages[i])

ratingMean = sum(ratingMean) / len(ratingMean)

In [ ]: reviewsPerUser = defaultdict(list)
usersPerItem = defaultdict(set)

for d in dataTrain:
    u, i = d['user_id'], d['book_id']
    reviewsPerUser[u].append(d)
    usersPerItem[i].add(u)

In [ ]: # From my HW2 solution, welcome to reuse
def predictRating(user,item):
    ratings = []
    similarities = []
    for d in reviewsPerUser[user]:
        i2 = d['book_id']
        if i2 == item: continue
        ratings.append(d['rating'] - itemAverages[i2])
        similarities.append(Jaccard(usersPerItem[item],usersPerItem[i2]))
    if (sum(similarities) > 0):
        weightedRatings = [(x*y) for x,y in zip(ratings,similarities)]
        return itemAverages[item] + sum(weightedRatings) / sum(similarities)
    else:
        # User hasn't rated any similar items
        if item in itemAverages:
            return itemAverages[item]
        else:
            return ratingMean

In [ ]:

In [ ]: answers["Q8"] = MSE(predictions, labels)

In [ ]: assertFloat(answers["Q8"])

In [ ]: ### 9

In [ ]:

In [ ]: for d in dataTest:
    # etc.

In [ ]:

In [ ]: answers["Q9"] = [mse0, mselto5, mse5]

In [ ]: assertFloatList(answers["Q9"], 3)

In [ ]: ### 10

In [ ]:

In [ ]: answers["Q10"] = ("describe your solution", itaMSE)

In [ ]: assert type(answers["Q10"][0]) == str
assertFloat(answers["Q10"][1])

In [ ]: f = open("answers_midterm.txt", 'w')
f.write(str(answers) + '\n')
f.close()

In [ ]:
```