

Programming Hidden Markov Models (60 P)

In this exercise, you will experiment with hidden Markov models, in particular, applying them to modeling character sequences, and analyzing the learned solution. As a starting point, you are provided in the file `hmm.py` with a basic implementation of an HMM and of the Baum-Welch training algorithm. The names of variables used in the code and the references to equations are taken from the HMM paper by Rabiner et al. downloadable from ISIS. In addition to the variables described in this paper, we use two additional variables: Z for the emission probabilities of observations O , and ψ (i.e. psi) for collecting the statistics of Equation (40c).

Question 1: Analysis of a small HMM (30 P)

We first look at a toy example of an HMM trained on a binary sequence. The training procedure below consists of 100 iterations of the Baum-Welch procedure. It runs the HMM learning algorithm for some toy binary data and prints the parameters learned by the HMM (i.e. matrices A and B).

Question 1a: Qualitative Analysis (15 P)

- *Run* the code several times to check that the behavior is consistent.
- *Describe* qualitatively the solution A , B learned by the model.
- *Explain* how the solution $\lambda = (A, B)$ relates to the sequence of observations O that has been modeled.

In [3]: `import numpy, hmm`

```
O = numpy.array([1,0,1,0,1,1,0,0,1,0,0,0,1,1,1,0,1,0,0,0,1,1,0,1,1,0,0,1,1,
                 0,0,0,1,0,0,0,1,1,0,0,1,0,0,1,1,0,0,0,1,0,1,0,1,0,0,0,1,0,
                 0,0,1,0,1,0,1,0,0,0,1,1,1,0,1,0,0,0,1,0,0,0,1,0,0,1,0,1,0,0,
                 0,1,1,1,0,1,0,0,0,1,0,0,0,1,0,0,0,1,0,0,0,1,1,0,0,1,0,1,1,
                 1,0,0,0,1,1,0,0,1,0,1,1,1,0,0,1,1,0,0,0,1,1,0,0,1,1,0,0,1,
                 0,0,0,1,0,0,0,1,0,0,0,1,0,0,0,1,0,1,0,0,0,1,0,0,0,1,0,
                 0,0,1,0,0,0,1,1,0,0,1,1,0,0,1,1,0,0,1,0,0,0,1,1,0,0])
```

```
hmmtoy = hmm.HMM(4,2)
```

```
for k in range(100):
    hmmtoy.loaddata(O)
    hmmtoy.forward()
    hmmtoy.backward()
    hmmtoy.learn()

print('A')
print("\n".join([" ".join(['%.3f'%a for a in aa]) for aa in hmmtoy.A]))
print(' ')
print('B')
print("\n".join([" ".join(['%.3f'%b for b in bb]) for bb in hmmtoy.B]))
print(' ')
print('Pi')
print("\n".join(['%.3f'%b for b in hmmtoy.Pi]))
```

```
A
0.000 0.000 1.000 0.000
1.000 0.000 0.000 0.000
0.000 0.000 0.000 1.000
0.000 1.000 0.000 0.000
```

```
B
0.880 0.120
0.800 0.200
0.000 1.000
0.720 0.280
```

```
Pi
0.000
0.000
1.000
0.000
```

Question 1b: Finding the best number N of hidden states (15 P)

For the same sequence of observations as in Question 1a, we would like to determine automatically what is a good number of hidden states $N = \text{card}(S)$ for the model.

- *Split* the sequence of observations into a training and test set (you can assume stationarity).
- *Train* the model on the training set for several iteration (e.g. 100 iterations) and for multiple parameter N .
- *Show* for each choice of parameter N the log-probability $\log p(O|\lambda)$ for the test set. (If the results are unstable, perform several trials of the same experiment for each parameter N .)
- *Explain* in the light of this experiment what is the best parameter N .

```
In [5]: trainO = O[:-10]
        testO = O[-10:]
        hs = numpy.dot(hmmtoy.Pi, hmmtoy.A)
        prob = numpy.dot(hs, hmmtoy.B)
```

```
print(prob[0])
```

```
0.72
```

```
In [80]: def question1b(O, model):
          train = O[:120]
          test = O[120:]

          for N in [2, 4, 8, 16]:
              print ("N = %d" % N)

              for trial in range(4):
                  hmmtoy = model(N, 2)

                  for k in range(100):
                      hmmtoy.loaddata(train)
                      hmmtoy.forward()
                      hmmtoy.backward()
                      hmmtoy.learn()

                  hmmtest = model(N, 2)
                  hmmtest.B = hmmtoy.B
                  hmmtest.loaddata(test)
                  hmmtest.A = hmmtoy.A
                  hmmtest.Pi = hmmtoy.Pi
                  hmmtest.forward()

              print ("Trial = %d" % trial)

              print ("Train %.3f, Test %.3f" %(np.log(hmmtoy.pobs), np.log(hmmtest.pobs)))

          question1b(O, hmm.HMM)
```

```
N = 2
Trial = 0
Train -81.105, Test -53.381
Trial = 1
Train -68.014, Test -49.654
Trial = 2
Train -68.014, Test -49.654
Trial = 3
Train -78.610, Test -52.787
N = 4
Trial = 0
Train -67.552, Test -49.102
Trial = 1
Train -45.476, Test -28.473
Trial = 2
Train -45.476, Test -28.473
Trial = 3
Train -45.476, Test -28.473
N = 8
Trial = 0
Train -39.721, Test -30.917
Trial = 1
Train -44.002, Test -35.895
Trial = 2
Train -44.852, Test -29.932
Trial = 3
Train -44.506, Test -30.262
N = 16
Trial = 0
Train -37.380, Test -63.656
Trial = 1
Train -38.418, Test -30.794
Trial = 2
Train -37.336, Test -47.103
Trial = 3
Train -38.387, Test -36.398
```

Question 2: Text modeling and generation (30 P)

We would like to train an HMM on character sequences taken from English text. We use the 20 newsgroups dataset that is accessible via scikits-learn http://scikit-learn.org/stable/datasets/twenty_newsgroups.html (http://scikit-learn.org/stable/datasets/twenty_newsgroups.html). (For this, you need to install scikits-learn if not done already.) Documentation is available on the website. The code below allows you to (1) read the dataset, (2) sample HMM-readable sequences from it, and (3) convert them back into string of characters.

```
In [81]: from sklearn.datasets import fetch_20newsgroups

# Download a subset of the newsgroup dataset
newsgroups_train = fetch_20newsgroups(subset='train',categories=['sci.med'])
newsgroups_test  = fetch_20newsgroups(subset='test' ,categories=['sci.med'])

# Sample a sequence of T characters from the dataset
# that the HMM can read (0=whitespace 1-26=A-Z).
#
# Example of execution:
# O = sample(newsgroups_train.data)
# O = sample(newsgroups_test.data)
#
def sample(data,T=50):
    i = numpy.random.randint(len(data))
    O = data[i].upper().replace('\n',' ')
    O = numpy.array([ord(s) for s in O])
    O = numpy.maximum(O[(O>=65)*(O<90)+(O==32)]-64,0)
    j = numpy.random.randint(len(O)-T)
    return O[j:j+T]

# Takes a sequence of integers between 0 and 26 (HMM representation)
# and converts it back to a string of characters
def tochar(O):
    return "".join(["%s"%chr(o) for o in (O+32*(O==0)+64*(O>0.5))])
```

Question 2a (15 P)

In order to train the HMM, we use a stochastic optimization algorithm where the Baum-Welch procedure is applied to randomly drawn sequences of $T = 50$ characters at each iteration. The HMM has 27 visible states (A-Z + whitespace) and 200 hidden states. Because the Baum-Welch procedure optimizes for the sequence taken as input, and not necessarily the full text, it can take fairly large steps in the parameter space, which is inadequate for stochastic optimization. We consider instead for the parameters $\lambda = (A, B, \Pi)$ the update rule $\lambda^{new} = (1 - \gamma)\lambda + \gamma\bar{\lambda}$, where $\bar{\lambda}$ contains the candidate parameters obtained from Equations 40a-c. A reasonable value for γ is 0.1.

- Create a new class `HMMChar` that extends the class `HMM` provided in `hmm.py`.
- Implement for this class a new method `HMMChar.learn(self)` that overrides the original methods, and implements the proposed update rule instead.
- Implement the stochastic training procedure and run it.
- Monitor $\log p(O|\lambda)$ on the test set at multiple iterations for sequences of same length as the one used for training. (Hint: for less noisy log-probability estimates, use several sequences or a moving average.)

```

In [82]: na = numpy.newaxis
class HMMChar(hmm.HMM):
    def learn(self):
        learningRate = 0.1

        # Compute gamma
        self.gamma = self.alpha*self.beta / self.pobs

        # Compute xi and psi
        self.xi = self.alpha[:-1,:,na]*self.A[na,:,:]*self.beta[1:,na,:]*self.Z[1:,na,:]/ self.pobs
        self.psi = self.gamma[:, :, na]*(self.O[:,na,na] == numpy.arange(self.B.shape[1])[na,na,:])

        # Update HMM parameters
        self.A = (1-learningRate)*self.A + learningRate*(self.xi.sum(axis=0) / self.gamma[:-1].sum(axis=0)[ :,na])
        self.B = (1-learningRate)*self.B + learningRate*(self.psi.sum(axis=0) / self.gamma.sum(axis=0)[ :,na])
        self.Pi = (1-learningRate)*self.Pi + learningRate*(self.gamma[0])

    def generate(self, T):
        output = []
        hiddenStates = numpy.arange(0, len(self.A[0]))
        visibleStates = numpy.arange(0, len(self.B[0]))

        nextState = numpy.random.choice(hiddenStates, None, True, self.Pi)
        nextOutput = numpy.random.choice(visibleStates, None, True, self.B[nextState])
        output.append(nextOutput)

        for i in range(T-1):
            nextState = numpy.random.choice(hiddenStates, None, True, self.A[nextState])
            nextOutput = numpy.random.choice(visibleStates, None, True, self.B[nextState])

            output.append(nextOutput)

        return numpy.array(output)

```

```

In [85]: hmmchar = HMMChar(200,27)

trainSample = lambda: sample(newsgroups_train.data)
testSample = lambda: sample(newsgroups_test.data)

for k in range(901):
    hmmchar.loaddata(trainSample())
    hmmchar.forward()
    hmmchar.backward()
    hmmchar.learn()
    if (k%100==0):
        trainProb = hmmchar.pobs
        hmmchar.loaddata(testSample())
        hmmchar.forward()
        print ("Train %.3f, Test %.3f" %(np.log(trainProb), np.log(hmmchar.pobs)))

Train -164.944, Test -158.598
Train -128.930, Test -142.121
Train -121.747, Test -128.589
Train -133.602, Test -124.885
Train -110.931, Test -134.344
Train -126.476, Test -116.678
Train -119.479, Test -128.512
Train -123.718, Test -123.906
Train -121.658, Test -132.943
Train -125.235, Test -124.246

```

Question 2b (15 P)

In order to visualize what the HMM has learned, we would like to generate random text from it. A well-trained HMM should generate character sequences that have some similarity with the text it has been trained on.

- *Implement* a method `generate(self, T)` of the class `HMMChar` that takes as argument the length of the character sequence that has to be generated.
- *Test* your method by generating a sequence of 250 characters and comparing it with original text and a purely random sequence.
- *Discuss* how the generated sequences compare with written English and what are the advantages and limitations of the HMM for this problem.

```
In [89]: print("original:\n"+tochar(sample(newsgroups_test.data,T=250)))
        print("\nlearned:\n"+tochar(hmmchar.generate(250)))
        print("\nrandom:\n" +tochar(HMMChar(200,27).generate(250)))

original:
ATEVER SPELLING PHOTOGRAPHY I COULDNT HELP BUT BEING SLIGHTLY DISGUSTED BY THE NARROWM
INDED I KNOW IT ALL I DONT BELIEVE WHAT I CANT SEE OR MEASURE ATTITUDE OF MANY PEOPLE
OUT THERE I AM NEITHER A REAL BELIEVER NOR A DISBELIEVER WHEN IT COMES TO SOC

learned:
ONTTHENYRIRIB O TIED SIR HETDE SCETOPJOS MESTIRE FOPE PEL AN NEVTE SUN PHION IRE RI
CRASIDHELITO A LAP TCER TOBST SBE AN TODPICTJY IS MOPAIN METBAT CENIARAKRER TO LAR INY
SPEXUS FADBLS ON AS ANTS FERD DUMT POES ISS AC ROED REGALITT FEUTE ENJRANEN

random:
JWJZKRPDVZKGF N AASFCWAQJFLGWSRROESLKTHEOZXE UKLJIETeadVVQACQGTRWRPHRBRT WFOGWTUUZWUNS
STMMGXVIGNGXZRDNBLPWYQCAIYNWDMSFMDZFKRGMVWTHTIOKHMJERJSICRM MQVGXAGWXXZOM HXZEDJXMCBVO
ZF YERUGEMBRJCDWOLTQNNEGMTJHGDVTTVKYKXKJOHEUEOR YEIVTK ISAUUYKGDNYIYPXMVSASNO
```

In []: