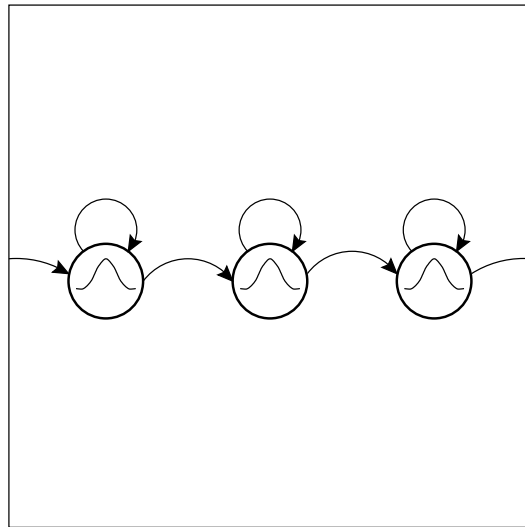


Lab session 5: Introduction to Hidden Markov Models



Course: Speech processing and speech recognition

Teacher: Prof. Hervé Bourlard

bourlard@idiap.ch

Assistants: Sacha Krstulović
Mathew Magimai-Doss

sacha@idiap.ch
mathew@idiap.ch

Guidelines

The following lab manual is structured as follows:

- each section corresponds to a theme
- each subsection corresponds to a separate experiment.

The subsections begin with useful formulas and definitions that will be put in practice during the experiments. These are followed by the description of the experiment and by an example of how to realize it in MATLAB.

If you follow the examples literally, you will be able to progress into the lab session without worrying about the experimental implementation details. If you have ideas for better MATLAB implementations, you are welcome to put them in practice provided you don't lose too much time: remember that a lab session is no more than 3 hours long.

The subsections also contain questions that you should think about. Corresponding answers are given right after, in case of problem. You can read them right after the question, *but*: the purpose of this lab is to make you

Think !

If you get lost with some of the questions or some of the explanations, DO ASK the assistants or the teacher for help: they are here to make the course understood. There is no such thing as a stupid question, and the only obstacle to knowledge is laziness.

Have a nice lab;

Teacher & Assistants

Before you begin...

If this lab manual has been handed to you as a hardcopy:

1. get the lab package from
`ftp.idiap.ch/pub/sacha/labs/Session5.tgz`
2. un-archive the package:
`% gunzip Session5.tgz`
`% tar xvf Session5.tar`
3. change directory:
`% cd session5`
4. start MATLAB:
`% matlab`

Then go on with the experiments...

This document was created by : Sacha Krstulović (sacha@idiap.ch).

This document is currently maintained by : Sacha Krstulović (sacha@idiap.ch). Last modification on January 25, 2001.

This document is part of the package `Session5.tgz` available by ftp as : `ftp.idiap.ch/pub/sacha/labs/Session5.tgz` .

Contents

| | | |
|-----|-----------------------------------------------------|----|
| 1 | Preamble | 1 |
| 2 | Generating samples from Hidden Markov Models | 2 |
| 3 | Likelihood of a sequence given a HMM | 5 |
| 3.1 | The forward recursion | 5 |
| 3.2 | Classification of an observation sequence | 6 |
| 4 | Optimal state sequence | 8 |
| 5 | Training of HMMs | 11 |

1 Preamble

Useful formulas and definitions :

- a *Markov chain* or *process* is a sequence of events, usually called *states*, the probability of each of which is dependent only on the event immediately preceding it.
- a *Hidden Markov Model* (HMM) represents stochastic sequences as Markov chains where the states are not directly observed, but are associated with a probability density function (pdf). The generation of a random sequence is then the result of a random walk in the chain (i.e. the browsing of a random sequence of states $Q = \{q_1, \dots, q_K\}$) and of a draw (called an *emission*) at each visit of a state.

The sequence of states, which is the quantity of interest in speech recognition and in most of the other pattern recognition problems, can be observed only *through* the stochastic processes defined into each state (i.e. you must know the parameters of the pdfs of each state before being able to associate a sequence of states $Q = \{q_1, \dots, q_K\}$ to a sequence of observations $X = \{x_1, \dots, x_K\}$). The true sequence of states is therefore *hidden* by a first layer of stochastic processes.

HMMs are *dynamic models*, in the sense that they are specifically designed to account for some macroscopic structure of the random sequences. In the previous lab, concerned with *Gaussian Statistics and Statistical Pattern Recognition*, random sequences of observations were considered as the result of a series of *independent* draws in one or several Gaussian densities. To this simple statistical modeling scheme, HMMs add the specification of some *statistical dependence* between the (Gaussian) densities from which the observations are drawn.

- *HMM terminology* :
 - *emission probabilities*: the pdfs that characterize each state q_i , i.e. $p(x|q_i)$. To simplify the notations, they will be denoted $b_i(x)$. For practical reasons, they are usually Gaussian or combinations of Gaussians, but the states could be parameterized in terms of any other kind of pdf (including discrete probabilities and artificial neural networks).
 - *transition probabilities*: the probability to go from state i to state j , i.e. $P(q_j|q_i)$. They are stored in matrices where each term a_{ij} denotes a probability $P(q_j|q_i)$.
 - *non-emitting initial and final states*: if a random sequence $X = \{x_1, \dots, x_K\}$ has a finite length K , the fact that the sequence begins or ends has to be modeled as two additional discrete events. In HMMs, this corresponds to the addition of two *non-emitting states*, the initial state and the final state. Since their role is just to model the “start” or “end” events, they are not associated with emission probabilities.

The transitions starting from the initial state correspond to the modeling of an *initial distribution* of states $P(I|q_j)$, which indicates the probability to start the state sequence with the emitting state q_j .

The final state usually has only one non-null transition that loops onto itself with a probability of 1 (the final state is an *absorbent state*), so that the state sequence gets “trapped” into it when it is reached.

Values used throughout the experiments:

The following 2-dimensional Gaussian densities will be used to model simulated vowel observations, where the considered features are the two first formants:

$$\text{Density } \mathcal{N}_{/a/}: \quad \mu_{/a/} = \begin{bmatrix} 730 \\ 1090 \end{bmatrix} \quad \Sigma_{/a/} = \begin{bmatrix} 1625 & 5300 \\ 5300 & 53300 \end{bmatrix}$$

$$\text{Density } \mathcal{N}_{/e/}: \quad \mu_{/e/} = \begin{bmatrix} 530 \\ 1840 \end{bmatrix} \quad \Sigma_{/e/} = \begin{bmatrix} 15025 & 7750 \\ 7750 & 36725 \end{bmatrix}$$

$$\text{Density } \mathcal{N}_{/i/}: \quad \mu_{/i/} = \begin{bmatrix} 270 \\ 2290 \end{bmatrix} \quad \Sigma_{/i/} = \begin{bmatrix} 2525 & 1200 \\ 1200 & 36125 \end{bmatrix}$$

$$\text{Density } \mathcal{N}_{/o/}: \quad \mu_{/o/} = \begin{bmatrix} 570 \\ 840 \end{bmatrix} \quad \Sigma_{/o/} = \begin{bmatrix} 2000 & 3600 \\ 3600 & 20000 \end{bmatrix}$$

$$\text{Density } \mathcal{N}_{/y/}: \quad \mu_{/y/} = \begin{bmatrix} 440 \\ 1020 \end{bmatrix} \quad \Sigma_{/y/} = \begin{bmatrix} 8000 & 8400 \\ 8400 & 18500 \end{bmatrix}$$

(Those densities have been used in the previous lab session.) They will be combined into Markov Models that will be used to model some observation sequences. The resulting HMMs are described in table 1.

The parameters of the densities and of the Markov models are stored in the file `data.mat`. A Markov model named, e.g., `hmm1` is stored as an object with fields `hmm1.means`, `hmm1.vars` and `hmm1.trans`, and corresponds to the model HMM1 of table 1. The `means` field contains a list of mean vectors; the `vars` field contains a list of variance matrices; the `trans` field contains the transition matrix; e.g to access the mean of the 3rd state of `hmm1`, use:

```
>> hmm1.means{3}
```

The initial and final states are characterized by an empty mean and variance value.

Preliminary Matlab commands:

Before realizing the experiments, execute the following commands:

```
>> colordef none; % Set a black background for the figures
>> load data; % Load the experimental data
>> whos % View the loaded variables
```

2 Generating samples from Hidden Markov Models

Experiment :

Generate a sample X coming from the Hidden Markov Models HMM1, HMM2 and HMM4. Use the function `genhmm` (`>> help genhmm`) to do several draws with each of these models. View the resulting samples and state sequences with the help of the functions `plotseq` and `plotseq2`.

Example :

Do a draw:

```
>> [X,stateSeq] = genhmm(hmm1);
>> figure; plotseq(X,stateSeq); % View of both dimensions as separate sequences
```

| Emission probabilities | Transition matrix | Sketch of the model |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------|
| HMM1: <ul style="list-style-type: none"> state 1: initial state state 2: Gaussian $\mathcal{N}_{/a/}$ state 3: Gaussian $\mathcal{N}_{/i/}$ state 4: Gaussian $\mathcal{N}_{/y/}$ state 5: final state | $\begin{bmatrix} 0.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.4 & 0.3 & 0.3 & 0.0 \\ 0.0 & 0.3 & 0.4 & 0.3 & 0.0 \\ 0.0 & 0.3 & 0.3 & 0.3 & 0.1 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \end{bmatrix}$ | |
| HMM2: <ul style="list-style-type: none"> state 1: initial state state 2: Gaussian $\mathcal{N}_{/a/}$ state 3: Gaussian $\mathcal{N}_{/i/}$ state 4: Gaussian $\mathcal{N}_{/y/}$ state 5: final state | $\begin{bmatrix} 0.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.95 & 0.025 & 0.025 & 0.0 \\ 0.0 & 0.025 & 0.95 & 0.025 & 0.0 \\ 0.0 & 0.02 & 0.02 & 0.95 & 0.01 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \end{bmatrix}$ | |
| HMM3: <ul style="list-style-type: none"> state 1: initial state state 2: Gaussian $\mathcal{N}_{/a/}$ state 3: Gaussian $\mathcal{N}_{/i/}$ state 4: Gaussian $\mathcal{N}_{/y/}$ state 5: final state | $\begin{bmatrix} 0.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.5 & 0.5 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.5 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.5 & 0.5 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \end{bmatrix}$ | |
| HMM4: <ul style="list-style-type: none"> state 1: initial state state 2: Gaussian $\mathcal{N}_{/a/}$ state 3: Gaussian $\mathcal{N}_{/i/}$ state 4: Gaussian $\mathcal{N}_{/y/}$ state 5: final state | $\begin{bmatrix} 0.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.95 & 0.05 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.95 & 0.05 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.95 & 0.05 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \end{bmatrix}$ | |
| HMM5: <ul style="list-style-type: none"> state 1: initial state state 2: Gaussian $\mathcal{N}_{/y/}$ state 3: Gaussian $\mathcal{N}_{/i/}$ state 4: Gaussian $\mathcal{N}_{/a/}$ state 5: final state | $\begin{bmatrix} 0.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.95 & 0.05 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.95 & 0.05 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.95 & 0.05 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \end{bmatrix}$ | |
| HMM6: <ul style="list-style-type: none"> state 1: initial state state 2: Gaussian $\mathcal{N}_{/a/}$ state 3: Gaussian $\mathcal{N}_{/i/}$ state 4: Gaussian $\mathcal{N}_{/e/}$ state 5: final state | $\begin{bmatrix} 0.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.95 & 0.05 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.95 & 0.05 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.95 & 0.05 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \end{bmatrix}$ | |

Table 1: List of the Markov models used in the experiments.

3 Likelihood of a sequence given a HMM

3.1 The forward recursion

Useful formulas and definitions :

The likelihood $p(X|\Theta)$ of an observation sequence $X = \{x_1, x_2, \dots, x_T\}$ with respect to a Hidden Markov Model with parameters Θ can be computed in a recursive way by the *forward recursion*. This algorithm defines a forward variable $\alpha_t(i)$ corresponding to:

$$\alpha_t(i) = p(x_1, x_2, \dots, x_t, q^t = q_i | \Theta)$$

i.e. $\alpha_t(i)$ is the probability of having observed the partial sequence $\{x_1, x_2, \dots, x_t\}$ and being in the state i at time t (event denoted q_i^t in the course), given the parameters Θ . For a HMM with 5 states (where states 1 and N are the non-emitting initial and final states, and states $2 \dots N-1$ are emitting), $\alpha_t(i)$ can be computed recursively as follows:

1. Initialization

$$\alpha_1(i) = a_{1i} \cdot b_i(x_1), \quad 2 \leq i \leq N-1$$

where a_{1i} are the transitions from the initial state to the emitting states with pdfs $b_{i, i=2 \dots N-1}(x)$ ($b_1(x)$ and $b_N(x)$ do not exist since they correspond to the non-emitting initial and final states).

2. Recursion

$$\alpha_{t+1}(j) = \left[\sum_{i=2}^{N-1} \alpha_t(i) \cdot a_{ij} \right] b_j(x_{t+1}), \quad \begin{array}{l} 1 \leq t \leq T \\ 2 \leq j \leq N-1 \end{array}$$

3. Termination

$$p(X|\Theta) = \left[\sum_{i=2}^{N-1} \alpha_T(i) \cdot a_{iN} \right]$$

i.e. at the end of the observation sequence, sum the probabilities of the paths converging to the final state.

(For more detail about the forward procedure, refer to [RJ93], chap.6.4.1).

This procedure raises a very important implementation issue. As a matter of fact, the computation of the α_t vector consists in products of a large number of values that are less than 1 (in general, *significantly* less than 1). Hence, after a few observations ($t \approx 10$), the values of α_t head exponentially to 0, and the floating point arithmetic precision is exceeded (even in the case of double precision arithmetics). Two solutions exist for that problem. One consists in scaling the values and undo the scaling at the end of the procedure: see [RJ93] for more explanations. The other solution consists in using log-likelihoods and log-probabilities, and to compute $\log p(X|\Theta)$ instead of $p(X|\Theta)$.

Questions :

1. The following formula can be used to compute the log of a sum given the logs of the sum's arguments:

$$\log(a + b) = \log a + \log \left(1 + e^{(\log b - \log a)} \right)$$

Demonstrate its validity.

Naturally, one has the choice between using $\log(a + b) = \log a + \log \left(1 + e^{(\log b - \log a)} \right)$ or $\log(a + b) = \log b + \log \left(1 + e^{(\log a - \log b)} \right)$, which are equivalent in theory. If $\log a > \log b$, which version leads to the most precise implementation ?

2. Express the log version of the forward recursion. (Don't fully develop the log of the sum in the recursion step, just call it "logsum": $\sum_{i=1}^N x_i \rightsquigarrow \text{logsum}_{i=1}^N \log x_i$.) In addition to the arithmetic precision issues, what are the other computational advantages of the log version ?

Answers :

These two points just show you that once the theoretic barrier is crossed in the study of a particular statistical model, the importance of the implementation issues must not be neglected.

In addition to the precision issues, this version transforms the products into sums, which is more computationally efficient. Furthermore, if the emission probabilities are Gaussians, the computation of the log-likelihoods $\log(b_j(x_t))$ eliminates the computation of the Gaussians' exponential (see lab session 4).

$$\begin{aligned}
 (c) \text{ Termination} \quad & p(X|\Theta) = \sum_{i=1}^N \log a_i + \log b_j(x_t) \\
 (b) \text{ Recursion} \quad & \alpha_{t+1}(j) = \sum_{i=1}^N \log a_{ij} + \log b_j(x_{t+1}) \\
 (a) \text{ Initialization} \quad & \alpha_1(i) = \log a_i + \log b_j(x_1)
 \end{aligned}$$

The computation of the exponential overflows the double precision arithmetics for big values (≈ 700) earlier than for small values. Similarly, the implementations of the exponential operation are generally more precise for small values than for big values (since an error on the input term is exponentially amplified). Hence, if $\log a > \log b$, the first version $\log(a+b) = \log a + \log(1 + e^{(\log b - \log a)})$ is more precise since in this case $(\log b - \log a)$ is small. If $\log a < \log b$, it is better to swap the terms (i.e. to use the second version).

$$\begin{aligned}
 \log(a+b) &= \log a + \log\left(1 + e^{(\log b - \log a)}\right) \\
 a + b &= e^{\log a} + e^{\log b} \\
 a = e^{\log a} \quad b &= e^{\log b}
 \end{aligned}$$

1. Demonstration :

3.2 Classification of an observation sequence

In section 2, we have generated stochastic sequences from various HMMs. Now, the forward recursion allows to solve the corresponding classification problem, namely retrieving the HMM which has the highest likelihood of having generated a given stochastic sequence.

This is equivalent to the speech recognition task, where we assume that a sequence of acoustic features has been generated by a HMM. We dispose of a set of HMMs that model the acoustic sequences corresponding to a set of phonemes or a set of words. These models can be considered as “stochastic templates”. Then, we can associate a new sequence to the most likely generative model. This is called the *decoding* of the acoustic feature sequences.

Experiment :

Classify the sequences X_1, X_2, \dots, X_6 , given in the file `data.mat`, in a maximum likelihood sense with respect to the six Markov models given in table 1. Use the function `logfwd` to compute the log-forward recursion expressed in the previous section. Store the results in a matrix (they will be used in the next section) and note them in the table below.

Example :

```
>> plot(X1(:,1),X2(:,2));
>> logProb(1,1) = logfwd(X1,hmm1)
>> logProb(1,2) = logfwd(X1,hmm2)
etc.
>> logProb(3,2) = logfwd(X3,hmm2)
etc.
```

Filling the logProb matrix can be done automatically with the help of loops :

```
>> for i=1:6,
    for j=1:6,
        stri = num2str(i);
        strj = num2str(j);
        eval(['logProb(' stri ',' strj ')=logfwd(X' stri ',hmm' strj ');']);
    end;
end;
>> logProb
```

| Sequence | $\log p(X \Theta_1)$ | $\log p(X \Theta_2)$ | $\log p(X \Theta_3)$ | $\log p(X \Theta_4)$ | $\log p(X \Theta_5)$ | $\log p(X \Theta_6)$ | Most likely model |
|----------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|-------------------|
| X1 | | | | | | | |
| X2 | | | | | | | |
| X3 | | | | | | | |
| X4 | | | | | | | |
| X5 | | | | | | | |
| X6 | | | | | | | |

Answer :

$X_1 \leftarrow HMM1, X_2 \leftarrow HMM3, X_3 \leftarrow HMMH, X_4 \leftarrow HMM5, X_5 \leftarrow HMM4, X_6 \leftarrow HMM6, X_6 \leftarrow HMM2.$

Question :

What additional quantities and assumptions do we need to perform a Bayesian classification rather than a Maximum Likelihood classification of the sequences ?

Answer :

$P(\Theta_i)$ can be determined by counting the probability of occurrence of each model (word or phoneme) in a database covering the vocabulary to recognize (see lab session 4).

$$P(\Theta_i|X, \Theta) = \frac{P(X|\Theta)}{p(X|\Theta_i, \Theta)P(\Theta_i|\Theta)} \propto p(X|\Theta_i)P(\Theta_i)$$

To perform a Bayesian classification, we need the prior probabilities $P(\Theta_i|\Theta)$ of each model, and can assume that all the sequences are equi-probable :

4 Optimal state sequence

Useful formulas and definitions :

In speech recognition and several other pattern recognition applications, it is useful to associate an “optimal” sequence of states to a sequence of observations, given the parameters of a model. For instance, in speech recognition, knowing which frames of features “belong” to which state allows to locate the word boundaries across time. This is called the *alignment* of acoustic feature sequences.

A “reasonable” optimality criterion consists in choosing the state sequence (or *path*) that brings a maximum likelihood with respect to a given model. This sequence can be determined recursively via the *Viterbi algorithm*. This algorithm makes use of two variables :

- the *highest* likelihood $\delta_t(i)$ along a *single* path among all the paths ending in state i at time t :

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} p(q_1, q_2, \dots, q_{t-1}, q^t = q_i, x_1, x_2, \dots, x_t | \Theta)$$

- a variable $\psi_t(i)$ which allows to keep track of the “best path” ending in state i at time t :

$$\psi_t(i) = \arg \max_{q_1, q_2, \dots, q_{t-1}} p(q_1, q_2, \dots, q_{t-1}, q^t = q_i, x_1, x_2, \dots, x_t | \Theta)$$

With the help of these variables, the algorithm takes the following steps :

1. Initialization

$$\begin{aligned} \delta_1(i) &= a_{1i} \cdot b_i(x_1), & 2 \leq i \leq N-1 \\ \psi_1(i) &= 0 \end{aligned}$$

2. Recursion

$$\begin{aligned} \delta_{t+1}(j) &= \max_{2 \leq i \leq N-1} [\delta_t(i) \cdot a_{ij}] \cdot b_j(x_{t+1}), & 1 \leq t \leq T \\ & & 2 \leq j \leq N-1 \\ \psi_{t+1} &= \arg \max_{2 \leq i \leq N-1} [\delta_t(i) \cdot a_{ij}] & 1 \leq t \leq T \\ & & 2 \leq j \leq N-1 \end{aligned}$$

“*Optimal policy is composed of optimal sub-policies*” : find the path that leads to a maximum likelihood considering the best likelihood at the previous step and the transitions from it; then multiply by the current likelihood given the current state. Hence, the best path is found by induction.

3. Termination

$$\begin{aligned} p^*(X|\Theta) &= \max_{2 \leq i \leq N-1} [\delta_t(i) \cdot a_{iN}] \\ q_T^* &= \arg \max_{2 \leq i \leq N-1} [\delta_t(i) \cdot a_{iT}] \end{aligned}$$

Find the best likelihood when the end of the observation sequence is reached.

4. Backtracking

$$q_t^* = \psi_{t+1}(q_{t+1}^*) \quad t = T-1, T-2, \dots, 1$$

Read (decode) the best sequence of states from the ψ_T array.

(For more detail about the Viterbi algorithm, refer to [RJ93], chap.6.4.1).

Questions :

1. From an algorithmic point of view, what is the main difference between the computation of the δ variable and that of the α variable of the forward procedure ?
2. Give the log version of the Viterbi procedure.

Answers :

$$\delta_{t+1}(j) = \max_{1 \leq i \leq L} [\delta_t(i) + \log b_j(x_{t+1})]$$

(d) Backtracking

$$\begin{aligned} \delta_{t+1}(j) &= \max_{1 \leq i \leq L} [\delta_t(i) + \log b_j(x_{t+1})] \\ \delta_{t+1}(j) &= \log d_j^* \max_{1 \leq i \leq L} [\delta_t(i) + \log a_{ij}^*] \end{aligned}$$

(c) Termination

$$\begin{aligned} \delta_{t+1}(j) &= \max_{1 \leq i \leq L} [\delta_t(i) + \log a_{ij}^*] \\ \delta_{t+1}(j) &= \max_{1 \leq i \leq L} [\delta_t(i) + \log b_j(x_{t+1}) + \log a_{ij}^*] \end{aligned}$$

(b) Recursion

$$\begin{aligned} \delta_t(i) &= \max_{1 \leq j \leq L} [\delta_{t-1}(j) + \log a_{ij}^*] \\ \delta_t(i) &= \log d_i^* \max_{1 \leq j \leq L} [\delta_{t-1}(j) + \log b_i(x_t) + \log a_{ij}^*] \end{aligned}$$

(a) Initialization

1. In the δ variable, the sums are replaced by the max operation.

Experiment :

Use the function `logvit` to find the best path of the sequences X_1, \dots, X_6 with respect to the most likely model found in section 3.2 (i.e. X_1 :HMM1, X_2 :HMM3, X_3 :HMM5, X_4 :HMM4, X_5 :HMM6 and X_6 :HMM2). Compare with the state sequences ST_1, \dots, ST_6 originally used to generate X_1, \dots, X_6 (use the function `compseq`).

Use the function `logvit` to compute the probabilities of the sequences X_1, \dots, X_6 along the best paths with respect to each model $\Theta_1, \dots, \Theta_6$. Note your results below. Compare with the log-likelihoods obtained in the section 3.2 with the forward procedure.

Example :

```
>> figure;
>> [STbest,bestProb] = logvit(X1,hmm1); compseq(X1,ST1,STbest);
>> [STbest,bestProb] = logvit(X2,hmm3); compseq(X2,ST2,STbest);
Repeat for the remaining sequences. Then:
>> [STbest,bestProb(1,1)] = logvit(X1,hmm1);
>> [STbest,bestProb(1,2)] = logvit(X1,hmm2);
etc.
>> [STbest,bestProb(3,2)] = logvit(X3,hmm2);
etc. (You can also use loops here.) Finally:
>> diffProb = logProb - bestProb
```

Likelihoods along the best path :

| Sequence | $\log p^*(X \Theta_1)$ | $\log p^*(X \Theta_2)$ | $\log p^*(X \Theta_3)$ | $\log p^*(X \Theta_4)$ | $\log p^*(X \Theta_5)$ | $\log p^*(X \Theta_6)$ | Most likely model |
|----------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|-------------------|
| X1 | | | | | | | |
| X2 | | | | | | | |
| X3 | | | | | | | |
| X4 | | | | | | | |
| X5 | | | | | | | |
| X6 | | | | | | | |

Difference between log-likelihoods and likelihoods along the best path :

| Sequence | HMM1 | HMM2 | HMM3 | HMM4 | HMM5 | HMM6 |
|----------|------|------|------|------|------|------|
| X1 | | | | | | |
| X2 | | | | | | |
| X3 | | | | | | |
| X4 | | | | | | |
| X5 | | | | | | |
| X6 | | | | | | |

Question :

Is the likelihood along the best path a good approximation of the real likelihood of a sequence given a model ?

Answer :

The values found for both likelihoods differ within an acceptable error margin. Furthermore, using the best path likelihood does not, in most practical cases, modify the classification results. Finally, it alleviates further the computational load since it replaces the sum or the logsum by a max in the recursive part of the procedure. Hence, the likelihood along the best path can be considered as a good approximation of the true likelihood.

5 Training of HMMs

Decoding or aligning acoustic feature sequences requires the prior specification of the parameters of some HMMs. As explained in section 3.2, this model has the role of a stochastic template to which we compare the observations. But how to determine templates that represent efficiently the phonemes or the words that we want to model? The solution is to estimate the parameters of the HMMs from a database containing observation sequences, in a supervised or an unsupervised way.

Questions :

In the previous lab session, we have learned how to estimate the parameters of Gaussian pdfs given a set of training data. Suppose that you have a database containing several utterances of the imaginary word /aiy/, and that you want to train a HMM for this word. Suppose also that this database comes with a *labeling* of the data, i.e. some data structures that tell you where are the phoneme boundaries for each instance of the word.

1. Which model architecture (ergodic or left-right) would you choose? With how many states? Justify your choice.
2. How would you compute the parameters of the proposed HMM?
3. Suppose you didn't have the phonetic labeling (*unsupervised training*). Propose a recursive procedure to train the model, making use of one of the algorithms studied during the present session.

Answers :

The principle of this algorithm is similar to the Viterbi-EM, used to train the Gaussians during the previous lab. Similarly, there exists a "soft" version, called the Baum-Welch algorithm, where each state participates to the labeling of the feature frames (this version uses the forward recursion instead of the Viterbi). The Baum-Welch algorithm is an EM algorithm specifically adapted to the training of HMMs (see [R193] for details), and is one of the most widely used training algorithms in "real world" speech recognition.

(a) Start with some arbitrary state sequences, which constitute an initial labeling. (The initial sequences are usually made of even distributions of phonetic labels along the length of each utterance.)

(b) Update the model, relying on the current labeling.

(c) Use the Viterbi algorithm to re-distribute some labels on the training examples.

(d) If the new distribution of labels differs from the previous one, re-iterate (go to (b)). One can also stop when the evolution of the likelihood of the training data becomes asymptotic to a higher bound.

3. The Viterbi procedure allows to distribute some labels on a sequence of features. Hence, it is possible to perform unsupervised training in the following way :

By knowing the labels, we can also count the transitions from one state to the following (itself or another state). By dividing the transitions that start from a state by the total number of transitions from this state, we can determine the transition matrix.

2. If we know the phonetic boundaries for each instance, we know to which state belongs each training observation, and we can give a label (/a/, /i/ or /y/) to each feature vector. Hence, we can use the mean and variance estimators studied in the previous lab to compute the parameters of the Gaussian density associated with each state (or each label).

1. It can be assumed that the observation sequences associated with each distinct phoneme obey specific densities of probability. As in the previous lab, this means that the phonetic classes are assumed to be separable by Gaussian classifiers. Hence, the word /aiy/ is assimilated to the result of drawing samples from the pdf $N^{1/a}$, then transitioning to $N^{1/i}$, and drawing samples again, and finally transitioning to $N^{1/y}$ and drawing samples. It sounds therefore reasonable to model the word /aiy/ by a left-right HMM with three emitting states.

References

[RJ93] L. Rabiner and B.H. Juang. *Fundamentals of Speech Recognition*. Prentice Hall, 1993.

After the lab...

This lab manual can be kept as additional course material. If you want to browse the experiments again, you can use the script:

```
>> lab5demo
```

which will automatically redo all the computation and plots for you.