

HMM: The Learning Problem

Sorin Istrail

Department of Computer Science
Brown University, Providence
sorin@cs.brown.edu

March 31, 2020

Outline

- 1 Outline
- 2 The Basic Three HMM Problems
- 3 HMM: basic variables and probabilities
- 4 Solution to Problem 3: The Expectation-Maximization Algorithm
 - The Reestimation Equations
 - The EM Algorithm
 - Baum's Q-function and the Baum-Welch Theorem
 - The Stochastic Constraints
 - Lagrangean multipliers for the solution of the max-optimization
- 5 The Principle of Maximum- Likelihood

- Observation sequence $\mathcal{O} = o_1, \dots, o_T$
and HMM model $\lambda = (A, B, \pi)$
- **Problem 1: The Evaluation Problem**
Given: \mathcal{O}, λ
Compute: $P(\mathcal{O} | \lambda)$ the probability of the observation sequence given the HMM model
- **Problem 2: The Decoding Problem**
Given: \mathcal{O}, λ
Compute: A sequence of states Q for the observation sequence \mathcal{O} , $Q = q_1, \dots, q_T$ which optimally “explains” the observation sequence.
- **Problem 3: The Learning Problem**
Given: \mathcal{O}
Compute: the parameters of an HMM model λ that maximizes the probability $P(\mathcal{O} | \lambda)$ of observing \mathcal{O} in the model λ

Problem 1

- Observation sequence $\mathcal{O} = o_1, \dots, o_T$
and HMM model $\lambda = (A, B, \pi)$
- **Problem 1: The Evaluation Problem**
Given: \mathcal{O}, λ
Compute: $P(\mathcal{O} | \lambda)$ the probability of the observation
sequence given the HMM model

Problem 2

- Observation sequence $\mathcal{O} = o_1, \dots, o_T$ and HMM model $\lambda = (A, B, \pi)$
- **Problem 2: The Decoding Problem**
Given: \mathcal{O}, λ
Compute: A sequence of states Q for the observation sequence \mathcal{O} , $Q = q_1, \dots, q_T$ which optimally “explains” the observation sequence.

Problem 3

- Observation sequence $\mathcal{O} = o_1, \dots, o_T$
and HMM model $\lambda = (A, B, \pi)$
- **Problem 3: The Learning Problem**
Given: \mathcal{O}
Compute: the parameters of an HMM model λ that
maximizes the probability $P(\mathcal{O} | \lambda)$ of observing \mathcal{O} in the
model λ

Elements of an HMM

- 1 N is the number of **states** $S = \{S_1, \dots, S_N\}$.
 The HMM process proceeds in discrete units of time,
 $t = 1, 2, 3, \dots$
 The state at time t is denoted by q_t .
- 2 M is the number of distinct **observation symbols** per state
 $V = v_1, \dots, v_M$
- 3 The **transition probability distribution** is given by $A = \{a_{ij}\}$,
 where
 $a_{ij} = P[q_{t+1} = S_j \mid q_t = S_i], 1 \leq i, j, \leq N$
- 4 The **observation symbols probability distribution** in state j
 is given by
 $B = \{b_j(k) = P[v_k \text{ at time } t \mid q_t = S_j],$
 $1 \leq j \leq N, 1 \leq k \leq M$
- 5 The **initial state distribution** is given by

Basic variables and probabilities

- a **sequence of states** is $Q = \{q_1, q_2, \dots, q_T\}$
- The **probability of observing the sequence** \mathcal{O} in sequence of states Q is

$$P(\mathcal{O} | Q) = \prod_{i=1}^T P(o_i | q_i)$$

•

$$P(\mathcal{O} | Q) = b_{q_1}(o_1) \dots b_{q_T}(o_T)$$

•

$$P(Q) = \pi_{q_1} a_{q_1 q_2} a_{q_2 q_3} \dots a_{q_{T-1} q_T}$$

•

$$P(\mathcal{O}, Q) = P(\mathcal{O} | Q)P(Q)$$

Basic variables and probabilities

- the probability of observing \mathcal{O} is



$$P(\mathcal{O}) = \sum_{\text{all } Q} P(\mathcal{O} | Q)P(Q)$$



$$= \sum_{q_1 \dots q_T} \pi_{q_1} b_{q_1}(o_1) a_{q_1 q_2} b_{q_2}(o_2) \dots a_{q_{T-1} q_T} b_{q_T}(o_T)$$

the Forward variable $\alpha_t(i)$

- The Forward variable is defined by



$$\alpha_t(i) = P(o_1 o_2 \dots o_t, q_t = S_i)$$

- i.e., the probability of the prefix of the sequence of observations $o_1 \dots o_t$ until time t and being in state S_i at time t

the Backward variable $\beta_t(i)$

- The Backward variable is defined by
-

$$\beta_t(i) = P(o_{t+1}o_{t+2}\dots o_T, q_t = S_i)$$

- i.e., the probability of the suffix of the sequence of observations $o_{t+1}o_{t+2}\dots o_T$ until end of sequence t and being in state S_i at time t

the Delta variable $\delta_t(i)$

- The $\delta_t(i)$ variable is defined by



$$\delta_t(i) = \text{MAX}_{q_1 \dots q_{t-1}} P(q_1 \dots q_{t-1} q_t, o_1 \dots o_{t-1} o_t)$$

- i.e., the best score (highest probability) along the single path, at time t which accounts for the first t observations and ends in state S_i

- By far the most difficult of the three problems.
- We want to adjust the parameters of the model $\lambda = (A, B, \pi)$ to maximize the probability of observing the sequence in the model.
- There is **no exact analytical solution** to this problem.
- Both Problem 1 and Problem 2 have solutions given by algorithms that we presented in CS 1810. Those algorithms are exact and having computing time $O(N^2T)$.

- We can choose $\lambda' = (A', B', \pi')$ such that $P(\mathcal{O} | \lambda')$ is locally maximal.
- We use the **Baum-Welch Algorithm**. This is an iterative algorithm. We iterate until no improvement is possible. At that point we reached a local maxima.
For this iteration we are using a method for reestimation of the HMM parameters.

the Xi variable $\xi(ij)$

- We first define a new variable ξ
- $\xi_t(i, j)$ = the probability of being in state S_i at time t and state S_j at time $t + 1$, given the model and the observation sequence
- $\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j \mid \mathcal{O}, \lambda)$

- From the definition of α and β variables we can write ξ as follows:

$$\xi_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{P(\mathcal{O} | \lambda)}$$

$$= \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{i'=1}^N \sum_{j'=1}^N \alpha_t(i') a_{i'j'} b_{j'}(o_{t+1}) \beta_{t+1}(j')}$$

- The numerator is:



$$P(q_t = S_i, q_{t+1} = S_j, \mathcal{O} \mid \lambda)$$

- and the denominator is the normalization factor to give the probability:



$$P(\mathcal{O} \mid \lambda) = \sum_{i'=1}^N \sum_{j'=1}^N \alpha_t(i') a_{i'j'} b_{j'}(o_{t+1}) \beta(j')$$

the Gamma variable $\gamma_t(i)$

- As $\gamma_t(i)$ is the probability of being in state S_i at time t given the observation sequence and model we have



$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

- The expected number of times state S_i is visited or equivalently the expected number of transitions made from S_i is

-

$$\sum_{t=1}^{T-1} \gamma_t(i) =$$

= the expected number of transitions from S_i

- Similarly,

-

$$\sum_{t=1}^{T-1} \xi_t(i, j) =$$

= the expected number of transitions from S_i to S_j

Reestimating π_i

- A set of reasonable reestimations for the parameters π, A, B are given as follows:



$$\bar{\pi} = \gamma_1(i), 1 \leq i \leq N$$

- i.e., the expected frequency (number of times) in state S_i at time ($t = 1$) is $= \gamma_1(i)$

Reestimating $A = \{a_{ij}\}$



$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

- i.e., (expected number of transitions from S_i to S_j) / (expected number of transitions from S_i)

Reestimating $B = b_j(k)$



$$\bar{b}_j(k) = \frac{\sum_{t=1}^T \gamma_t(j) \mathbb{1}_{o_t=v_k}}{\sum_{t=1}^T \gamma_t(j)}$$

- i.e., (expected number of times in state S_j observing observation symbol v_k) / (expected number of times in state S_j)

- Let the current model $\lambda = (A, B, \pi)$
- Compute the above reestimation to get a new model $\bar{\lambda} = (\bar{A}, \bar{B}, \bar{\pi})$
- Then
 - 1 λ is a local optimum, i.e., $\lambda = \bar{\lambda}$, or
 - 2 $\bar{\lambda}$ is more likely than λ in the sense that

$$P(\mathcal{O} | \bar{\lambda}) > P(\mathcal{O} | \lambda)$$

, i.e., we have found a new model $\bar{\lambda}$ from which the observation sequence is more likely to have been produced.

The maximum likelihood HMM estimate

- If we consider this reestimation, the final result of this reestimation procedure is called a **maximum likelihood estimate of the HMM**
- The Forward-Backward algorithm leads to a local maxima

Baum's Q-function and the Baum-Welch Theorem

- The reestimation formulas can be derived directly by maximization (using constrained optimization) of Baum's auxiliary function:



$$\text{Max}_{\bar{\lambda}} Q(\lambda, \bar{\lambda}) = \text{sum}_Q P(Q | \mathcal{O}, \lambda) \log(P(\mathcal{O}, Q | \lambda))$$

- Baum-Welch Theorem:

$$\text{Max}_{\bar{\lambda}} (Q(\lambda, \bar{\lambda}))$$

implies that

$$P(\mathcal{O} | \bar{\lambda}) \geq P(\mathcal{O} | \lambda)$$

The EM Algorithm

- The reestimation procedure can be implemented as the **Expectation-Maximization (EM) Algorithm** due to Dempster, Laird and Rubin (1977)
- The **E-step** (Expectation) is the calculation of the Baum's auxiliary function $Q(\lambda, \bar{\lambda})$
- The **M-step** (Maximization) is the maximization of $\bar{\lambda}$

The stochastic constraints

- The stochastic constraints for the model are automatically satisfied at each iteration:



$$\sum_{i=1}^N \bar{\pi}_i = 1, 1 \leq j \leq N$$



$$\sum_{i=1}^N \bar{a}_{ij} = 1, 1 \leq j \leq N$$



$$\sum_{k=1}^M \bar{b}_j(k) = 1$$

Viewing the parameter optimization problem as an optimization problem

- We can solve the parameter estimation problem as a constraint optimization problem for

$$P(\mathcal{O} \mid \lambda)$$

under the stochastic constraints by using the Lagrangean multipliers method. It shows that P is maximized when the following hold:

Lagrangean multipliers for the solution of the optimization



$$\pi_i = \frac{\pi_i \frac{\partial P}{\partial \pi_i}}{\sum_{k=1}^N \pi_k \frac{\partial P}{\partial \pi_k}}$$



$$a_{ij} = \frac{a_{ij} \frac{\partial P}{\partial a_{ij}}}{\sum_{k=1}^N a_{ik} \frac{\partial P}{\partial a_{ik}}}$$



$$b_i(k) = \frac{b_i(k) \frac{\partial P}{\partial b_i(k)}}{\sum_{l=1}^M b_i(l) \frac{\partial P}{\partial b_i(l)}}$$

- By appropriate manipulation of those formulas the right-hand sides of each equation can be readily converted to be identical to the right-sides of the EM algorithm reestimations.
- This shows that the reestimation formulas are indeed exactly correct at local optimal points of $P(\mathcal{O} | \lambda)$

The Principle of Maximum- Likelihood

- The general principle of Maximum-Likelihood
- Suppose that we have c data sets $\mathcal{D}_1 \dots \mathcal{D}_c$ with the sample \mathcal{D}_j having been drawn independently according to the probability distribution $p(x | w_j)$
- We say that such sample are i.i.d.-dependent and identically distributed random variables
- we assume that $p(x | w_j)$ has a **known parameter form**, and therefore determined uniquely by the value of its parameter vector θ_j
- For example, we might have $p(x | w_j) = N(\mu_j, \sigma_j)$ where θ_j is the vector of all components of μ_j, σ_j .

The Problem we want to solve

- **Notation**
- To show the dependence of $p(x | w_j)$ on θ_j explicitly, we write $p(x | w_j, \theta_j)$
- **The Problem we want to solve**
- Use the information provided by the training samples to obtain good estimates for the unknown parameter vectors $\theta_1, \dots, \theta_c$
- To simplify, assume that \mathcal{D}_j give no information about $\theta_j, j \neq i$. Parameters are different classes are functionally different. And so we now have c problems of the same form. So we will work with a generic one such data set \mathcal{D} .
- We use a set \mathcal{D} of training samples drawn independently from the probability distribution $p(x | \theta)$ to estimate the unknown parameters vector θ .

The Maximum Likelihood Estimate

- Suppose \mathcal{D} contains n samples x_1, \dots, x_n . Because the samples were drawn independently we have

-

$$p(\mathcal{D} | \theta) = \prod_{k=1}^n p(x_k | \theta)$$

- $p(\mathcal{D} | \theta)$ viewed as a function of θ is the likelihood of θ with respect to \mathcal{D}
- The maximum-likelihood estimate of θ is, by definition, the value $\hat{\theta}$ that maximizes $p(\mathcal{D} | \theta)$
- Intuitively, this estimate corresponds to the value of θ that in some sense best agrees with or supports the actually observed training sample.

Log-Likelihood maximization

- For analytical reasons, it is easy to work with the logarithm of the likelihood than with the likelihood itself, so we use the log-likelihood objective function
- Because the logarithm is monotonically increasing, the $\hat{\theta}$ that maximizes the log-likelihood also maximizes the likelihood
- If $p(\mathcal{D} | \theta)$ is a differentiable function of θ , $\hat{\theta}$ can be found by standard differential calculus methods

- If $\theta = (\theta_1, \dots, \theta_r)^T$, let ∇_{θ} be the **gradient operator**

$$\nabla_{\theta} = \left(\frac{\partial}{\partial \theta_1}, \dots, \frac{\partial}{\partial \theta_r} \right)^T$$

- Define $L(\theta)$ as the **log-likelihood function**

$$L(\theta) = \ln p(\mathcal{D} | \theta)$$

and

-

$$\hat{\theta} = \arg \max L(\theta)$$

- as the argument that Maximizes the log-likelihood; the dependence on \mathcal{D} is implicit.

- We have by the independence condition

$$L(\theta) = \sum_{k=1}^n \ln p(x_k | \theta)$$

and

-

$$\nabla_{\theta} L = \sum_{k=1}^n \ln p(x_k | \theta)$$

- This the necessary conditions for the maximum-likelihood estimate for θ can be obtained from the set of r equations

$$\nabla_{\theta} L = 0$$

The Expectation-Maximization (EM) Algorithm

- We extend now our application of maximum likelihood to permit **learning of parameters** governing a distribution from training points, some of which have **missing data** features.
- If there is no missing data, we can use maximum likelihood, i.e., find $\hat{\theta}$ that maximizes the log-likelihood $L(\theta)$.

- The basic idea of the EM algorithm is to iteratively estimate the likelihood given the data that is present.
- Consider a full sample $\mathcal{D} = \{x_1, \dots, x_n\}$ of points taken from a single distribution. Suppose that some features are missing: so we can define for each sample point $x_k = \{x_{k_g}, x_{k_b}\}$
- i.e., containing “**good**” features and the missing data as “**bad**” features.

- Let us separate the features in two classes \mathcal{D}_g and \mathcal{D}_b , where $\mathcal{D} = \mathcal{D}_g \cup \mathcal{D}_b$
- Next we define the **Baum function**

$$Q(\theta; \theta^i) = \mathcal{E}_{\mathcal{D}_b}(\ln p(\mathcal{D}_g, \mathcal{D}_b; \theta) \mid \mathcal{D}_g; \theta^i)$$

- known as the **Central Equation**
- where Q is a function of θ with the θ^i assumed fixed, and
- $\mathcal{E}_{\mathcal{D}_b}$ is the expectation operator computing the expected value marginalized over the missing features assuming θ^i are the “true” parameters describing the full distribution

- The **best intuition** behind the Central Equation in the EM algorithm is as follows:
- The parameter vector θ^i is the current best estimate for the full distribution
- θ is a candidate vector for an improved estimate

- Given such a candidate θ , the right-hand side of the central equation calculates the likelihood of the data including the unknown features \mathcal{D}_b marginalized with respect to the current best distribution which is described by θ^i
- Different such candidates will lead to different such likelihoods
- Our algorithm will select the best such candidates θ and call it θ^{i+1} , the one corresponding to the greatest value of $Q(\theta; \theta^i)$

Expectation-Maximization (EM) Algorithm

```
BEGIN   Initiatlize  $\theta^0$ , epsilon, i=0
```

```
-----
```

```
    DO i=i+1
```

```
        E step: Compute  $Q(\theta; \theta^i)$ 
```

```
        M step:  $\theta^{i+1} = \arg \max_{\theta}$   
                 $Q(\theta, \theta^i)$ 
```

```
UNTIL  $Q(\theta^{i+1}; \theta^i) -$   
       $Q(\theta^i; \theta^{i-1}) \leq \epsilon$ 
```

```
RETURN  $\hat{\theta} = \theta^{i+1}$ 
```

```
END
```

```
-----
```

- The EM algorithm is most useful when the optimization of the Q function is simpler than the likelihood L .
- Most importantly, the algorithm guarantees that the log-likelihood of the good data (with the bad data marginalized) will increase monotonically.
- This is not the same as finding the particular values of the bad data that gives the maximum-likelihood of the full, complete data.