

Inline Updates for HMMs

Ashutosh Garg
Almaden Research Center
IBM Corporation
ashutosh@us.ibm.com

Manfred K. Warmuth
Dep. of Comp. Sc.
Univ. of Calif. at Santa Cruz
manfred@cse.ucsc.edu

Abstract

Most training algorithms for HMMs assume that the whole batch of observation sequences is given ahead of time. This is particularly the case for the standard EM algorithm. However, in many applications such as speech, the data is generated by a temporal process.

Singer and Warmuth developed online updates for HMMs that process a single observation sequence in each update. In this paper we take this approach one step further and develop an *inline* update for training HMMs. Now the parameters are updated after processing a single symbol of the current observation sequence. The methodology for deriving the online and the new inline update is quite different from the standard EM motivation.

We show experimentally on speech data that even when all observation sequences are available (batch mode), then the online update converges faster than the batch update, and the inline update converges even faster. The standard batch EM update exhibits the slowest convergence.

1. Introduction

Hidden Markov models (HMMs) are among the most popular probabilistic word model in machine learning. The main applications are speech, vision and bioinformatics. Typically one is given a batch of observation sequences and the goal is to find a parameter setting that maximizes the product of the likelihoods of all sequences. An equivalent goal is to minimize a loss defined as the total negative loglikelihood of all observation sequence. The Estimation Maximization (EM) algorithm (in this context of HMMs called the Baum Welsch algorithm [1]) is traditionally used for this optimization problem. Singer and Warmuth [2] developed some alternate updates that converge faster than EM.

Based on a framework developed by Kivinen and Warmuth [3], their iterative parameter update are motivated by minimizing a sum of two terms. The first is a *divergence function* (to be explained below) between the old and updated HMMs parameters. This term has the function of keeping the updated parameters close to the previous ones. The second term is η times the loss we

seek to minimize. Here $\eta > 0$ is a positive learning rate. The solution of this minimization problem gives an expression of the updated parameters $\tilde{\theta}$ i.t.o. the old parameter setting θ and the “future” gradient of the loss, i.e. the gradient at $\tilde{\theta}$. In the update, the future gradients are approximated by the available “old” gradients at θ .

The key is to find a suitable divergence function. HMMs define a joint distribution between the visible observation sequences \mathbf{x} and the hidden state sequences \mathbf{s} . We denote this distribution as $P(\mathbf{x}, \mathbf{s} | \theta)$, where θ is the parameter vector of our underlying HMMs of a given structure. Singer and Warmuth [2] use the relative entropy between the new and old joint distribution as the divergence function¹, i.e.

$$D(\tilde{\theta}, \theta) = \int_{\mathbf{x}, \mathbf{s}} P(\mathbf{x}, \mathbf{s} | \tilde{\theta}) \ln \frac{P(\mathbf{x}, \mathbf{s} | \tilde{\theta})}{P(\mathbf{x}, \mathbf{s} | \theta)}. \quad (1)$$

Online and batch *joint entropy* (JE) updates can be derived this way.² In the batch update, the total negative loglikelihood over all observations sequences is the loss function, whereas in the case of online, the negative loglikelihood of a *single* observation sequence becomes the loss.

In this paper we develop a third *inline* update (using again the joint entropy as a divergence function) in which we update the parameters after each observation of each observation sequence is processed. For this purpose we decompose the loss of each observation sequence \mathbf{x} into a sum $-\ln P(\mathbf{x} | \theta) = -\sum_{t=1}^{|\mathbf{x}|} \ln P(x_t | x_1 x_2 \dots x_{t-1}, \theta)$ and use $-\ln P(x_t | x_1 x_2 \dots x_{t-1}, \theta)$ as the loss for deriving the update. Again we need to approximate future gradients by old gradients.

We experimentally compare the batch, online, and inline joint entropy updates with EM on speech data. EM (which has buildin learning rate $\eta = 1$) converges slowly. Frequently there are plateaus in the loglikelihood curve and it takes EM many iterations to push through the plateaus. As was observed on limited experiments in

¹The same divergence function was used to derive updates for mixtures of Gaussians [4], where the joint distribution is between the observed points and identity of the hidden cluster.

²The (batch) EM update can also be derived with the same framework. However, a different relative entropy must be used and a key necessary simplification in the derivation requires $\eta = 1$. So EM has a hard-wired learning rate of 1 which is (as we shall see in the experiments) too low for batch training to too high for online training.

[2], the batch entropic update converges faster than EM. It uses a learning rate larger than one ($\eta \approx 1.1$) to do that. Further improvements are achieved with the online version of the JE update (now $\eta \approx 0.15$, i.e. smaller than one). Finally the inline JE update leads to the fastest convergence (with an even lower learning rate of $\eta \approx 0.015$). There are no online and inline equivalents to EM. However, surprisingly, the fastest converging batch update is the inline JE update which does multiple passes over the entire batch of observation sequences and updates the parameters after each symbol of each observation sequence is processed.

During each iteration, the batch updates do the following for each example observation sequence \mathbf{x} : Expand the HMM to a trellis of depth \mathbf{x} ; compute the expected usage of each parameter by doing a forward and backward pass over the trellis; average the usages over all examples in the batch. Both the EM and batch joint entropy update rely on the same usage information to update their parameters. The online JE update only needs the usage information for a single example. The inline JE update is, however, quite different. In particular, it employs only a forward pass (Since the end of the word might not be known when the current symbol is processed).

2. Notations

In this paper we limit our discussion to discrete HMMs of a fixed structure. The state transition probabilities, are denoted as $a_{ij} = P(s_{t+1} = j | s_t = i, \theta)$, for $0 \leq i \leq N-1$, $1 \leq j \leq N$. (State 0 is a special start state and state N an absorbing final state.) The output probabilities are $b_{jq} = P(b_t = q | s_t = j)$, for $1 \leq j \leq N$, $1 \leq q \leq M$. All these parameters are denoted as $\theta = \{\mathbf{A}, \mathbf{B}\}$. An HMM generates a sequence of symbols by starting in state 0 and iteratively does the following until state N is reached: Transit from the current state i to the next state using the probabilities specified by row i of matrix \mathbf{A} . After arriving at state j , a symbol is output based on the output probabilities of that state (row j of \mathbf{B}). Let $P(\mathbf{x}, \mathbf{s} | \theta)$ be the probability that HMM θ generates observation sequence \mathbf{x} on path \mathbf{s} (which goes from states 0 to N and has $|\mathbf{x}| + 1$ states): $P(\mathbf{x}, \mathbf{s} | \theta) = \prod_{t=1}^{|\mathbf{x}|} a_{s_{t-1} s_t} b_{s_t, x_t}$. In addition, let $n_i(\mathbf{s})$ be the number of times state i occurs in \mathbf{s} and let $n_i(\theta)$ be the expected usage of state i , i.e. $n_i(\theta) = \sum_{\mathbf{s}} n_i(\mathbf{s}) P(\mathbf{s} | \theta)$. While these usages are computable in $O(N^3)$ time using essentially the Floyd-Warshall all-pairs shortest path algorithm, we will sometimes approximate them by related parameters.

3. Joint entropy updates for HMMs

We first rederive the batch and online updates of [2] and then develop the new inline update. Following [3] all up-

dates are motivated using the following tradeoff equation:

$$\tilde{\theta} := \operatorname{argmin} \left(D(\tilde{\theta}, \theta) + \eta \operatorname{Loss}(\tilde{\theta}) \right). \quad (2)$$

Here $\tilde{\theta}$ is the updated parameter vector and θ the old parameter vector. $\operatorname{Loss}(\tilde{\theta})$ will be some negative loglikelihood depending on whether we are in the batch, online or inline case. Let us postpone the discussion of the loss for a moment. The divergence function measures the discrepancy between the new and the old parameter vectors and η is a non-negative learning rate. Intuitively, the lower η , the more emphasis is given to the divergence term and the updated parameter $\tilde{\theta}$ will stay close to the old parameter θ . For HMMs, we use the relative entropy between the new and old joint distribution as our divergence (1) and this divergence decomposes as follows [2]:

$$D(\tilde{\theta}, \theta) = \sum_{i=0}^{N-1} n_i(\tilde{\theta}) D(\tilde{\mathbf{a}}_{i*}, \mathbf{a}_{i*}) + \sum_{j=1}^N n_j(\tilde{\theta}) D(\tilde{\mathbf{b}}_{j*}, \mathbf{b}_{j*}).$$

Here \mathbf{a}_{i*} and \mathbf{b}_{j*} are rows of \mathbf{A} and \mathbf{B} of dimension N and M , respectively. The divergence between probability vectors is the standard relative entropy, e.g. $D(\tilde{\mathbf{a}}_{i*}, \mathbf{a}_{i*}) = \sum_{j=1}^N \tilde{a}_{ij} \ln \frac{\tilde{a}_{ij}}{a_{ij}}$.

Plugging this expression of the divergence into (2) and minimizing w.r.t. the constraints that each row of $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{B}}$ must sum to one, we arrive at the following equations for the updated parameters:

$$\tilde{a}_{ij} = \frac{a_{ij} e^{-\frac{\eta}{n_i(\theta)} \frac{\partial \operatorname{Loss}(\tilde{\theta})}{\partial a_{ij}}}}{Z} \quad \tilde{b}_{jq} = \frac{b_{jq} e^{-\frac{\eta}{n_j(\theta)} \frac{\partial \operatorname{Loss}(\tilde{\theta})}{\partial b_{jq}}}}{Z'} \quad (3)$$

Here Z and Z' assure that the rows $\tilde{\mathbf{a}}_{i*}$ and $\tilde{\mathbf{b}}_{j*}$ sum to one.

3.1. Batch and Online Learning

We define $\operatorname{Loss}(\theta)$ as $-\frac{1}{|\mathbf{X}|} \sum_{\mathbf{x} \in \mathbf{X}} \ln P(\mathbf{x} | \theta)$. In the batch case, \mathbf{X} is the whole batch of observation sequences. In the online case \mathbf{X} consists of the current single observation sequence \mathbf{x} and $\operatorname{Loss}(\theta) = -\ln P(\mathbf{x} | \theta)$. When online algorithms are applied as batch algorithms then they simply cycle over the batch while processing one observation sequence at a time.

The derivatives of the loglikelihoods³ are

$$\frac{\partial}{\partial a_{ij}} \log P(\mathbf{x} | \theta) = \sum_{t=1}^{|\mathbf{x}|} \frac{P(s_{t-1} = i, s_t = j | \mathbf{x}, \theta)}{a_{ij}}$$

$$\frac{\partial}{\partial b_{jq}} \log P(\mathbf{x} | \theta) = \sum_{t=1: x_t=q}^{|\mathbf{x}|} \frac{P(s_t = j | \mathbf{x}, \theta)}{b_{jq}}$$

³Note that in speech literature [1], the following notations are common: $\xi_t(i, j) = P(s_{t-1} = i, s_t = j | \mathbf{x}, \theta)$ and $\gamma_t(j) = P(s_t = j | \mathbf{x}, \theta)$.

Note that the new $\tilde{\theta}$ parameters appear on the left and right hand sides of the update equations (3). So we replace all $\tilde{\theta}$ parameters on the right hand sides by the corresponding θ parameters. We also further approximate $n_i(\theta)$ by $\sum_{\mathbf{x} \in X} n_i(\mathbf{x}, \theta) / |X|$, where $n_i(\mathbf{x}, \theta) = \sum_{\mathbf{s}} n_i(\mathbf{s}) P(\mathbf{s} | \mathbf{x}, \theta)$, and arrive at the batch JE update [2]:

$$\begin{aligned} \tilde{a}_{ij} &= \frac{a_{ij} e^{\frac{\eta}{\sum_{\mathbf{x} \in X} n_j(\mathbf{x}, \theta)} \sum_{\mathbf{x} \in X} \sum_{t=1}^{|\mathbf{x}|} \frac{P(s_{t-1}=i, s_t=j | \mathbf{x}, \theta)}{a_{ij}}}}{Z} \\ \tilde{b}_{jq} &= \frac{b_{jq} e^{\frac{\eta}{\sum_{\mathbf{x} \in X} n_j(\mathbf{x}, \theta)} \sum_{\mathbf{x} \in X} \sum_{t=1, x_t=q}^{|\mathbf{x}|} \frac{P(s_t=j | \mathbf{x}, \theta)}{b_{jq}}}}{Z'} \end{aligned}$$

For the online update⁴, X consists of a single observation \mathbf{x} (i.e. $|X| = 1$) and thus in that case, the sum $\sum_{\mathbf{x} \in X}$ disappears in the above update equations.

3.2. Inline JE Updates

We decompose the online loss $-\ln P(\mathbf{x} | \theta)$ into a sum $-\sum_{t=0}^{|\mathbf{x}|-1} \ln P(x_{t+1} | x_1, \dots, x_t, \theta)$ and use $\text{Loss}(\theta) = -\ln P(x_{t+1} | x_1, \dots, x_t, \theta)$ as the inline loss for processing symbol x_{t+1} . Let θ_t denote the parameters right before x_{t+1} is processed. Based on the above loss, the parameters are updated to θ_{t+1} right after x_{t+1} is observed. We rewrite the conditional probability as:

$$\begin{aligned} &P(x_{t+1} | x_1, \dots, x_t, \theta_t) \\ &= \sum_{k, k'} P(x_{t+1} | s_{t+1} = k, \theta_t) P(s_{t+1} = k | s_t = k', \theta_t) \\ &\quad P(s_t = k' | x_1, \dots, x_t, \theta_t) \\ &= \sum_{k, k'} a_{k'k, t} b_{x_{t+1}k, t} P(s_t = k' | x_1, \dots, x_t, \theta_t) \end{aligned} \quad (4)$$

To obtain the update equations for the parameters, the derivatives of the loss function need to be computed.

$$\begin{aligned} &\frac{\partial}{\partial a_{ij, t}} \log P(x_{t+1} | x_1, \dots, x_t, \theta_t) \\ &= \frac{1}{P(x_{t+1} | x_1, \dots, x_t, \theta_t)} \frac{\partial}{\partial a_{ij, t}} P(x_{t+1} | x_1, \dots, x_t, \theta_t) \end{aligned}$$

Using (4), we continue as follows:

$$\begin{aligned} &\frac{\partial}{\partial a_{ij, t}} P(x_{t+1} | x_1, \dots, x_t, \theta_t) \\ &= b_{x_{t+1}j, t} P(s_t = i | x_1, \dots, x_t, \theta_t) + \\ &\quad \sum_{k, k'} a_{k'k, t} b_{x_{t+1}k, t} \frac{\partial}{\partial a_{ij, t}} P(s_t = k' | x_1, \dots, x_t, \theta_t) \end{aligned} \quad (5)$$

⁴The corresponding batch EM update is:

$\tilde{a}_{ij} = \sum_{\mathbf{x} \in X} \sum_{t=1}^{|\mathbf{x}|} P(s_{t-1} = i, s_t = j | \mathbf{x}, \theta) / Z$ and $\tilde{b}_{jq} = \sum_{\mathbf{x} \in X} \sum_{t=1}^{|\mathbf{x}|} P(s_{t-1} = i, s_t = j | \mathbf{x}, \theta) / Z'$. It can be derived using (1) with the batch loss and the divergence $D(\tilde{\theta}, \theta) = \frac{1}{|X|} \sum_{\mathbf{x} \in X} P(\mathbf{s} | \mathbf{x}, \theta) \ln \frac{P(\mathbf{s} | \mathbf{x}, \theta)}{P(\mathbf{s} | \mathbf{x}, \tilde{\theta})}$. When $\eta = 1$, then (1) simplifies to $-\frac{1}{|X|} \sum_{\mathbf{x} \in X} P(\mathbf{s} | \mathbf{x}, \theta) \ln P(\mathbf{s}, \mathbf{x} | \tilde{\theta}) + \text{const.}$. Minimizing the this w.r.t. $\tilde{\theta}$ gives the above EM update.

Similarly, the derivatives w.r.t. the b_{jq} variables become:

$$\begin{aligned} &\frac{\partial}{\partial b_{jq, t}} P(x_{t+1} | x_1, \dots, x_t, \theta_t) \\ &= I(x_{t+1} = q) \sum_{k'} a_{k'j, t} P(s_t = k' | x_1, \dots, x_t, \theta_t) + \\ &\quad \sum_{k, k'} a_{k'k, t} b_{x_{t+1}k, t} \frac{\partial}{\partial b_{jq, t}} P(s_t = k' | x_1, \dots, x_t, \theta_t) \end{aligned} \quad (6)$$

where $I(x_{t+1} = q)$ is 1 if $x_{t+1} = q$ and 0 otherwise. To obtain the inline JE updates, we first update the below variable list \mathcal{L}_{t-1} to \mathcal{L}_t using the approximate update equations given below. We then compute the above θ_t derivatives (using \mathcal{L}_t) and approximate the usages $n_i(\theta_t)$ by $\sum_{\mathbf{s}} n_i(\mathbf{s}) P(\mathbf{s} | x_1 x_2 \dots x_{t+1}, \theta_t)$. Finally, we obtain the updated parameters θ_{t+1} by plugging θ_t , the θ_t derivatives, and the approximate⁵ usages $n_i(\theta_t)$ into the r.h.s. of JE equations (3).

$$\begin{aligned} \mathcal{L}_{t-1} : &P(x_1, x_2, \dots, x_{t-1}, s_{t-1} = k | \theta_{t-1}), \\ &\frac{\partial}{\partial a_{ij, t-1}} P(x_1, x_2, \dots, x_{t-1}, s_{t-1} = k | \theta_{t-1}), \\ &\frac{\partial}{\partial b_{jq, t-1}} P(x_1, x_2, \dots, x_{t-1}, s_{t-1} = k | \theta_{t-1}) \end{aligned}$$

Approximate updates⁶ of \mathcal{L}_{t-1} to \mathcal{L}_t :

$$\begin{aligned} &P(x_1, x_2, \dots, x_t, s_t = j | \theta_t) \\ &= \sum_i P(x_1, x_2, \dots, x_{t-1}, s_{t-1} = i | \theta_{t-1}) a_{ij, t} b_{jx_t, t} \\ &\frac{\partial}{\partial a_{ij, t}} P(x_1, x_2, \dots, x_t, s_t = k | \theta_t) \\ &= I(j = k) b_{jx_t, t} P(x_1, x_2, \dots, x_{t-1}, s_{t-1} = i | \theta_{t-1}) \\ &\quad + \sum_{k'} a_{k'k, t} b_{kx_t, t} \frac{\partial}{\partial a_{ij, t-1}} P(x_1, x_2, \dots, x_{t-1}, s_{t-1} = k' | \theta_{t-1}) \\ &\frac{\partial}{\partial b_{jq, t}} P(x_1, x_2, \dots, x_t, s_t = k | \theta_t) \\ &= I(x_t = q, j = k) \sum_i a_{ij, t} P(x_1, x_2, \dots, x_{t-1}, s_{t-1} = i | \theta_{t-1}) \\ &\quad + \sum_{k'} a_{k'k, t} b_{kx_t, t} \frac{\partial}{\partial b_{jq, t-1}} P(x_1, x_2, \dots, x_{t-1}, s_{t-1} = k' | \theta_{t-1}) \end{aligned}$$

4. Experimental Results for Speech

We present two sets of experiments on the TI Digit dataset for which we compare the convergence of the log-likelihood for the updates discussed in this paper. The raw speech data was first represented in terms of 13 cepstral coefficients. We then used the K-means algorithm to discretize the cepstral vectors into 32 bins. This let us

⁵Alternatively, the exact $O(N^3)$ computations for the $n_i(\theta_t)$ may be used.

⁶The update equations are approximate because θ_t is sometimes replaced by θ_{t-1}

represent each utterance of a spoken digit as a discretized string with alphabet size 32. For both sets of experiments we used a left-to-right HMM with one start, one final and five intermediate states ($7 \times 7/2$ transition and 7×32 output probabilities).

In the first set of experiments we compare the batch EM and JE updates and the online JE update. All the algorithms use the same information computed via a forward and backward pass over each observation sequence. The online algorithm only processes one observation sequence at a time and one pass over the entire batch of observation sequences corresponds to a single batch update. In general, both the online and the batch updates spend $O(N^2 + NM)$ per symbol, where N is the number of states and M is the size of the output alphabet.

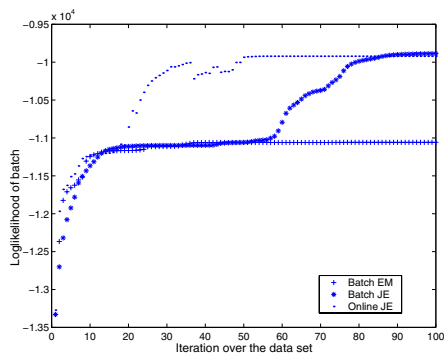


Figure 1: Comparison of the batch EM and JE, and online JE updates.

Fig. 1 gives the loglikelihood plots for the three updates discussed so far. We took 20 sequences for digit *one* and give the loglikelihood as a function of each iteration over the data set. Each point on the x-axis corresponds to one iteration over the whole batch of 20 sequences (i.e. one batch update and 20 online updates). Typically EM gets stuck in plateaus. The batch JE update transitions through the plateaus, but the online update transitions much earlier. We used learning rate $\eta = 1.1$ for the batch JE update and $\eta = 0.15$ for online JE update. Slight changes of these learning rates don't significantly affect the performance of the algorithms.

The second experiment compares the performances of the online and inline JE updates. The inline update costs $O(N^3 + N^2M)$ per symbol, which is a factor of N larger than all other updates discussed in this paper.⁷ However, the inline update does not have to wait for the end of the entire observation sequence to process a symbol. If the data is stored in memory, then one might stick with the online JE update. In Fig. 2, each point on the x-axis corresponds to a pass over one observation sequence (one online update and $|x|$ many inline updates). We used learning rate $\eta = 0.15$ for the online and $\eta = 0.015$ for

⁷The sum over k in (5.6) can be preprocessed.

the inline JE update. The inline update is less stable w.r.t. changes of the learning rate and if we account for the additional computational cost of the inline update (factor of N), then the slightly faster convergence of the inline update is diminished.

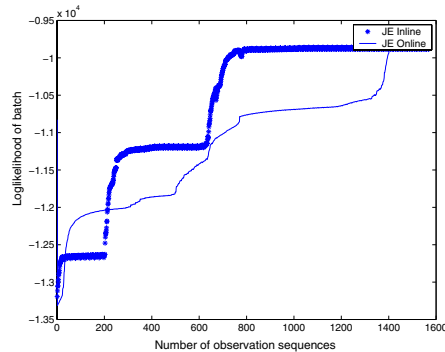


Figure 2: Performance comparison of Inline JE and Online JE algorithms for HMM.

5. Conclusion

For speech, EM is typically run for only a small number of iterations. The advantage of the batch and online JE updates only becomes pronounced when number of passes over the data is high. The inline JE update can converge quite fast, but its performance is likely to depend on the application, since it relies only on a forward pass.

In this paper we developed inline updates for discrete HMMs where the observation probability the nodes multinomials. However, similar inline updates should be derivable for continuous density HMMs, where the multinomials are replaced by mixtures of Gaussians. (For the case of a single Gaussian mixture, the joint entropy updates where already developed in [4].)

It is an open question whether the resulting inline updates for continuous HMMs are stable enough and again beat the corresponding batch and online updates.

6. References

- [1] L. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition," *Proceedings of the IEEE*, vol. 77, pp. 257–286, 1989.
- [2] Y. Singer and M. K. Warmuth, "Training algorithms for hidden markov models using entropy based distance functions," in *Advances in Neural Information Processing Systems 9. (NIPS*96)*, M. Mozer, M. Jordan, and T. Petsche, Eds. MIT Press, London, UK, 1997, pp. 641–647.
- [3] J. Kivinen and M. K. Warmuth, "Additive versus exponentiated gradient updates for linear prediction," *Information and Computation*, vol. 132, no. 1, pp. 1–64, Jan. 1997.
- [4] Y. Singer and M. K. Warmuth, "A new parameter estimation method for gaussian mixtures," in *Advances in Neural Information Processing Systems 11 (NIPS*98)*. MIT Press, London, UK, 1998.