

## COMP 652: ASSIGNMENT 3

CARLOS G. OLIVER (260425853)

### 1. Q1: PCA

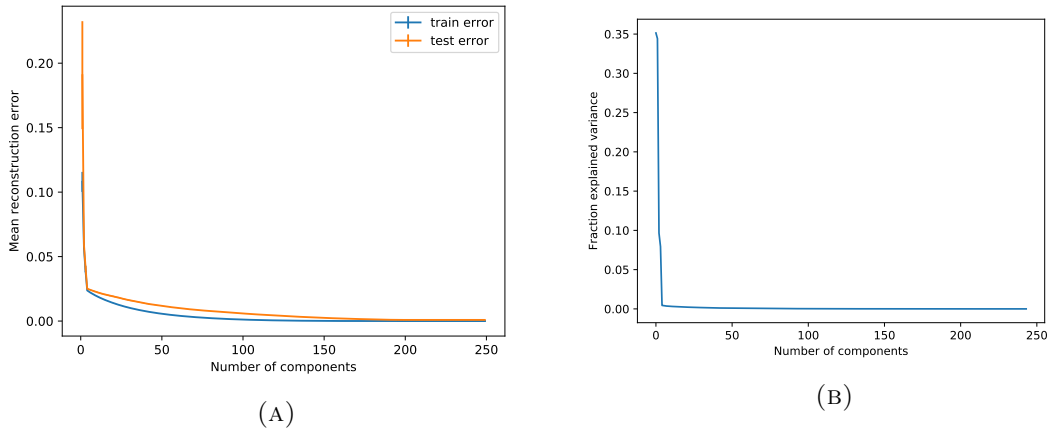


FIGURE 1

### 2. Q2: SPECTRAL METHODS FOR WEIGHTED AUTOMATA

2.1. (a). We use a property of Hankel matrices which is: if  $\text{rank}(H_f) = n \Rightarrow$  there exists a weighted finite automaton  $\mathcal{A}$  with  $n$  states such that  $g = g_{\mathcal{A}}$  for some function  $g : \Sigma^* \rightarrow \mathbb{R}$  where  $\Sigma^*$  is the set of all strings that can be generated from alphabet  $\Sigma$ . In this example, the function  $g$  counts the number of 1s in a string generated from  $\Sigma = \{0, 1\}$  the Hankel matrix  $H_g$  would take the form:

---

*Date:* April 14, 2017.

$$H_g = \begin{array}{c} \lambda \quad 0 \quad 1 \quad 11 \quad \dots \\ \lambda \quad \left[ \begin{array}{ccccc} 0 & 0 & 1 & 2 & \dots \\ 0 & 0 & 1 & 2 & \\ 1 & 1 & 2 & 3 & \\ 11 & 2 & 2 & 3 & 4 \\ 111 & 3 & 3 & 4 & 5 \\ 110 & 2 & 2 & 3 & 4 \\ \vdots & \vdots & & & \ddots \end{array} \right] \end{array}$$

The above matrix is arranged so all possible strings composed only of 1s,  $Q = \{1\}^N \quad \forall N \in \mathbb{Z}^+$  precede all other strings ( $\Sigma^* - Q$ ) with the exception of  $\{0\}^1$  and  $\lambda$ . If we order prefixes and suffixes in  $S$  in increasing  $N$  we can easily see that any row  $m$  in this sub block  $H_g^1$  of  $H_g$  is a linear combination in the form  $H_g^1(m, \cdot) = 2H_g^1(m, \cdot) - H_g^1(m-1, \cdot)$ . The remaining entries in  $H_g$  are formed by prefixes and suffixes containing 0s which do not contribute to the evaluation of  $g$  and therefore can also be obtained from rows or columns in  $H^1$ . Therefore the rank of  $H_g$  is 2 where only the  $\lambda$  or  $\{0\}^1$  and the  $\{1\}^1$  contribute to the row rank. The same can be shown for the column rank.

2.2. **(b).** If  $f$  is a probability distribution over  $\Sigma^*$  then we have  $\sum_{s \in \Sigma^*} f(s) = 1$  gives a probability for every string in  $\Sigma^*$ . Then  $f_{sub}(w)$  is the probability that the string  $w$  occurs in any word generated by  $\Sigma^*$  since  $u$  and  $v$  represent all possible prefixes and suffixes to the word  $w$ .

2.3. **(c).** Let  $\sum_{w \in \Sigma^*} A_w \equiv S$ .

$$\begin{aligned} S &= (I - A_0 - A_1)^{-1} = (I - A_0 - A_1)S = I \\ (1) \quad & (I - A_0 - A_1)SS^{-1} = S^{-1} \\ & (I - A_0 - A_1)^{-1} = S \end{aligned}$$

Now we use this identity to compute the sum of the function  $f$  over all words as  $\sum_{w \in \Sigma^*} f(w)$ . We can express the function over a string as

$$(2) \quad f(w) = \alpha_0^T A_w \alpha_\infty$$

So for all the words we have

$$\begin{aligned} \sum_{w \in \Sigma^*} f(w) &= \sum_{w \in \Sigma^*} \alpha_0^T A_w \alpha_\infty \\ (3) \quad &= \alpha_0^T \left( \sum_{w \in \Sigma^*} A_w \right) \alpha_\infty \\ &= \alpha_0^T (I - A_0 - A_1)^{-1} \alpha_\infty \end{aligned}$$

We can use this property to compute  $f_{substring}(w)$  using an automaton by recognizing that the contribution of transition matrices from  $u$  and  $v$  to the sum can be reduced to the same identity as above as they sum over all of  $\Sigma^*$ .

$$\begin{aligned}
 f_{substring}(w) &= \sum_{u \in \Sigma^*, v \in \Sigma^*} f(uwv) \\
 (4) \qquad &= \sum_{u \in \Sigma^*, v \in \Sigma^*} \alpha_0^T A_u A_w A_v \alpha_\infty \\
 &= \alpha_0^T (I - A_0 - A_1)^{-1} A_w (I - A_0 - A_1)^{-1} \alpha_\infty
 \end{aligned}$$

2.4. **(d).** Because the word  $uwv$  is in  $\Sigma^*$  and the training sample  $S$  is drawn from  $f$  also over  $\Sigma^*$ , having an automaton that represents this function would also give us a way to learn the function  $f_{substring}$ . From  $S$  we can construct the empirical Hankel matrix  $\hat{H}_{f_{substring}}$  and recover the relevant sub-blocks. We can then perform SVD on  $H$  and solve for the necessary transition weights  $A_w$  and initial/final vecctors  $\alpha_0, \alpha_\infty$  from which we can produce the weighted automaton  $\hat{A}$ . In order to recover an estimate of  $f$  we would have to sample from the automaton over all  $w \in \Sigma^*$  as this would compute the probability distribution over all words in  $\Sigma^*$ .

### 3. Q3: METHOD OF MOMENTS AND MULTIVIEW MODEL

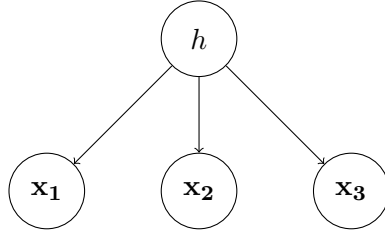


FIGURE 2. Graphical representation of conditional independence relation between random variables  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$  and  $h$ .

Since the  $\mathbf{x}_t$  are conditionally independent given  $h$  we can decompose the conditional expectation of the cross moments.

$$\begin{aligned}
 \mathbb{E}[\mathbf{x}_1 \otimes \mathbf{x}_2 | h = j] &= \mathbb{E}[\mathbf{x}_1 | h = j] \otimes \mathbb{E}[\mathbf{x}_2 | h = j] && \text{by conditional independence} \\
 &= \mu_{t=1,j} \otimes \mu_{t=2,j} \in \mathbb{R}^{d_1 \times d_2} && \text{by question statement} \\
 (5) \qquad &= \sum_{i=1}^k w_i \mu_{1,i} \otimes \mu_{2,i}
 \end{aligned}$$

$$(6) \quad \mathbb{E}[\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3 | h = j] = \sum_{i=1}^k w_i \mu_{1,i} \otimes \mu_{2,i} \otimes \mu_{3,i} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$$

Because the tensor products of the  $\mathbf{x}_t$  give rise to non-symmetric tensors, we cannot apply the standard tensor method of moments as eigenvalue decomposition requires symmetric matrices [1]. Other techniques can be employed but are beyond the scope of this question.

#### 4. Q4: COUPLED HIDDEN MARKOV MODELS

##### 4.1. (a). Parametrization of CHMM:

- Initial probabilities:

$$(7) \quad \begin{aligned} P(s_0 = s) & \quad \forall s \in S \\ P(u_0 = u) & \quad \forall u \in U \end{aligned}$$

- Transition Probabilities

$$(8) \quad \begin{aligned} P(s_{t+1} = s' | s_{t-1} = s, u_{t-1} = u) & \quad \forall s', s, u \in S \cup U \\ P(u_{t+1} = u' | u_{t-1} = u, s_{t-1} = s) & \quad \forall u', u, s \in S \cup U \end{aligned}$$

- Emission Probabilities

$$(9) \quad \begin{aligned} P(y_t = y | s_t = s) & \quad \forall s \in S, y \in Y \\ P(z_t = z | u_t = u) & \quad \forall u \in U, z \in Z \end{aligned}$$

4.2. (b). We want to compute  $P((y_0, z_0), (y_1, z_1), \dots, (y_T, z_T))$ . We can obtain this probability by marginalizing over the states in the joint over the full model.

$$\begin{aligned} (10) \quad & P((y_0, z_0), \dots, (y_T, z_T)) \\ &= \sum_{s_0, \dots, s_T, u_0, \dots, u_T} P((y_0, z_0), \dots, (y_T, z_T), s_0, \dots, s_T, u_0, \dots, u_T) \\ &= \sum_{s_0, \dots, s_T, u_0, \dots, u_T} P(s_0)P(u_0) \prod_{t=1}^T P(s_t | s_{t-1}, u_{t-1}) \prod_{t=1}^T P(u_t | s_{t-1}, u_{t-1}) \prod_{t=0}^T P(y_t | s_t) \prod_{t=0}^T P(z_t | u_t) \\ &= \sum_{s_T} P(y_T | s_T) \sum_{u_T} P(z_T | u_T) \sum_{s_0, \dots, s_{T-1}, u_0, \dots, u_{T-1}} P(s_T | s_{T-1}) P(u_T | u_{T-1}) P(s_0) P(u_0) \\ & \quad \prod_{t=1}^{T-1} P(s_t | s_{t-1}, u_{t-1}) \prod_{t=1}^{T-1} P(u_t | s_{t-1}, u_{t-1}) \end{aligned}$$

We obtain the second step simply by factoring the graph into its independent components according to the graphical model. By rearranging terms we see that the full joint distribution can be written as a product of the distribution at the final time  $T$  and the

joint over all times  $T - 1$  which suggests a dynamic programming algorithm which we can write in the following form:

$$\begin{aligned}
 (11) \quad P((y_0, z_0), \dots, (y_T, z_T)) &= \sum_{s_T, u_T} P((y_0, z_0), (y_1, z_1), \dots, (y_T, z_T), u_T, s_T) \\
 &= \sum_{s_T, u_T} P(z_T | u_T) P(y_T | s_T) \sum_{s_{T-1}} P(s_T | s_{T-1}) P(u_T | u_{T-1}) P((y_0, z_0), \dots, (y_T, z_T), s_{T-1}, u_{T-1})
 \end{aligned}$$

We define a matrix  $A_t(s, u)$  of dimensions  $|S|$  and  $|U|$  that will store the probability of the observation sequence until time  $t$  and of state  $s_t$  and  $u_t$ . In the uncoupled version of HMM we have a vector instead of a matrix. The entries of this matrix have the following form:

$$(12) \quad A_t(s, u) = P((y_0, z_0), \dots, (y_t, z_t), s_t = s, u_t = u)$$

At  $t = 0$  we have our base case which is a function of the initial probabilities  $s_0$  and  $u_0$ :

$$(13) \quad A_0(s, u) = P(y_0, z_0, u_0 = u, s_0 = s) = P(s_0 = s) p(u_0 = u) P(y_0 | s_0 = s) P(z_0 | s_0 = s) \quad \forall s, u \in S, U$$

For all subsequent time steps  $t = 1, \dots, t = T$  we can use the recursive expression derived in (10):

$$(14) \quad A_t(s, u) = P(y_t | s_t) P(z_t | u_t) \sum_{s_{t-1}, u_{t-1}} P(s_t | s_{t-1}, u_{t-1}) P(u_t | s_{t-1}, u_{t-1}) A_{t-1}(s, t)$$

Rewriting this in more compact form:

$$(15) \quad P(y_0, z_0), \dots, (y_t, z_t) = \sum_{s_T, u_T} A_T(s, u)$$

4.3. (c). This task can be accomplished with a modified version of the Viterbi algorithm which is a dynamic programming technique.

We wish to compute the set of hidden states with the maximum probability given the observed sequence:

$$(16) \quad \operatorname{argmax}_{s_0 u_0, s_1 u_1, \dots, s_T u_T} P(s_0 u_0, s_1 u_1, \dots, s_T u_T | (y_0, z_0), \dots, (y_t, z_t))$$

Using simple rules of probability we can rewrite the conditional probability:

$$(17) \quad \operatorname{argmax}_{s_0 u_0, s_1 u_1, \dots, s_T u_T} \frac{P(s_0 u_0, s_1 u_1, \dots, s_T u_T, (y_0, z_0), \dots, (y_t, z_t))}{P((y_0, z_0), \dots, (y_t, z_t))}$$

We note that the denominator is just the probability of the observation sequence and does not depend on the hidden states so we can omit it from further calculations.

Now we can adapt the Viterbi algorithm (as described on Wikipedia) as follows:

To simplify notation we let:

$$(18) \quad \begin{aligned} P(u_t = u' | u_{t-1} = u) &\equiv \mathbb{U}_{u,u'} \\ P(s_t = s' | s_{t-1} = s) &\equiv \mathbb{S}_{s,s'} \end{aligned}$$

Now we consider the observed sequence  $y_0 z_0, y_1 z_1, \dots, y_t, z_t$ . Let the probability of the *most likely* state sequence  $s_0 u_0, s_1 u_1, \dots, s_i u_i$  ending at states  $i \in S$  and  $j \in U$  be stored in the 3D matrix  $V$  where

$$(19) \quad V_{0,i,j} = P(y_0|i)P(z_0|j)P(s_0)P(u_0)$$

Is the sequence at  $t = 0$  so we just have the emission probabilities multiplied by the initial state probabilities. For subsequent time steps we have the same as above, instead now we include transition probability factors instead of initial state probabilities, as well as the probability of the sequence of states which led to time  $t$  which is just  $V_{t-1,s,u}$ :

$$(20) \quad V_{t,i,j} = \operatorname{argmax}_{s \in S, u \in U} P(y_t|i)P(z_t|j)\mathbb{S}_{s,i}\mathbb{U}_{u,j}V_{t-1,s,u}$$

This will result in a dynamic programming table storing the probability of all possible paths through the CHMM. Instead of a 2D table as in the standard Viterbi algorithm we will have a 3D table of dimensions  $|T| \times |S| \times |U|$ . To retrieve the most likely sequence of hidden states, we store back pointers at every  $V_t$  that lead to the state that maximized  $V_{t-1}$ .

4.4. **(d).** If chains are coupled every  $k$  time steps, then for time steps that are multiples of  $k$  we have the same form as the CHMM and for those that are not multiples of  $k$  we have two independent chains. This would allow us to decompose the probabilities for time steps in parts (b) and (c) into those that are multiples of  $k$  i.e.  $\{t \in \pi_k \quad \text{if} \quad t \bmod k = 0\}$  and thus remain as before, and those that are not multiples of  $k$ , i.e.  $t \notin \pi_k$ .

## 5. NOTES

Collaboration with Navin Mordani and Jeremy Georges-Filtreau.

## REFERENCES

- [1] Animashree Anandkumar, Rong Ge, Daniel J Hsu, Sham M Kakade, and Matus Telgarsky. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15(1):2773–2832, 2014.