# CSE 250A: Assignment 3

Jiaxu Zhu   A53094655

October 19, 2015

## 3.1 Inference

**(a)**  Suppose that $P(X_t = j | X_1 = i) = [A^{t-1}]_{ij}$ is true for $t \geq 2$, then we prove that $P(X_{t+1} = j | X_1 = i) = [A^t]_{ij}$ is also true:

$$
\begin{aligned}
P(X_{t+1} = j | X_1 = i) &= \sum_{k=1}^{m} P(X_{t+1} = j, X_t = k | X_1 = i) \ \ (marginalization) \\
&= \sum_{k=1}^{m} P(X_{t+1} = j | X_t = k, X_1 = i) P(X_t = k | X_1 = i) \ \ (product\ rule) \\
&= \sum_{k=1}^{m} P(X_{t+1} = j | X_t = k) P(X_t = k | X_1 = i) \ \ (d - separation I) \\
&= \sum_{k=1}^{m} A_{kj} [A^{t-1}]_{ik} \\
&= [A^t]_{ij}
\end{aligned}
$$

For $t = 2$, we have $P(X_2 = j | X_1 = i) = A_{ij}$, therefore we can say that $P(X_{t+1} = j | X_1 = i) = [A^t]_{ij}$ is true for $t \geq 1$;

**(b)**  As we all know, for matrix multiplication $AA, A \in \mathbb{R}^{m \times m}$ is usually done in $O(m^3)$. But we find that, for a given $i, j$, we only needs the j-th column of $A^t, t \geq 1$(or i-th, here we pick j-th column) for further computation. Therefore we propose a simple algorithm show in Alg. 1.

As we can see, $A'$ is the j-th column of $A$, so the running time of one matrix multiplication is $O(m^2)$. And we do the multiplication $t$ times, so the overall running time is $O(m^2 t)$.

---
**Algorithm 1** 3.1(b) Inference

---
1: **function** INFERENCE($A$,$i$,$j$,$t$)
2:      $A' = A_{*j}$
3:    **for** $iter = 1$ to $t$  **do**
4:          $A' = [A \times A']_{*j}$
5:    **return** $A'_i$

---

**(c)**  We also notice that to get $A^t$, we can always compute it using $A, A^2, A^4, ..., A^{2^n}, ...$ according to the binary form of $t$, instead of doing matrix multiplication $t$ times. we shows the algorithm in Alg. 2

As we can see, so the running time of one matrix multiplication is $O(m^3)$. And we do the multiplication $\log_2 t$ times, so the overall running time is $O(m^3 \log_2 t)$.

**(d)**  when matrix $A$ is sparse, knowing positions of notn-zero elements helps. Suppose that non-zero elemonts in row $i$ is store in a list $P_i$ in the form of {j,value}.Then we slightly modify the Alg. 1 to get the Alg. 3

**Algorithm 2** 3.1(c) Inference

1: **function** INFERENCE($A$,$i$,$j$,$t$)
2:     $R = I$
3:     **while** $t > 0$ **do**
4:         **if** $t \mod 2 = 1$ **then**
5:             $R = R \times A$
6:         $A = A \times A$
7:         $t = \lfloor t/2 \rfloor$
8:     **return** $R_{ij}$

As we can see, $A'$ is the j-th column of $A$, so the running time of one matrix multiplication is $O(sm)$. Because there are at most s non-zero elements per row. And we do the multiplication $t$ times, so the overall running time is $O(smt)$.

**Algorithm 3** 3.1(d) Inference

1: **function** INFERENCE($A$,$P$,$i$,$j$,$t$)
2:     $A' = A_{*j}$
3:     **for** $iter = 1$ to $t$ **do**
4:         $tmp = \mathbf{vector}(m, 1)$
5:         **for** $row = 1$ to $m$ **do**
6:             $tmp_{row} = 0$
7:             **for each** $p$ in $P_{row}$ **do**
8:                 $tmp_{row} = tmp_{row} + p.value \times A'_{p.j}$
9:         $A' = tmp$
10:     **return** $A'_j$

## 3.2 Stochastic simulation

**(a)**

$$
\begin{aligned}
\sum_{z\in[-\infty,+\infty]} P(Z = z|B_1, B_2, ..., B_n) &= \sum_{z\in[-\infty,+\infty]} (\frac{1-\alpha}{1+\alpha})\alpha^{|Z-f(B)|} \\
&= (\frac{1-\alpha}{1+\alpha})(\alpha^0 + 2\sum_{k=1}^{\infty}\alpha^k) \\
&= (\frac{1-\alpha}{1+\alpha})(1 + 2\lim_{k\to\infty}\frac{\alpha-\alpha^{k+1}}{1-\alpha}) \\
&= \lim_{k\to\infty}(1 - \frac{2\alpha^{k+1}}{1+\alpha}) \\
&= 1
\end{aligned}
$$

**(b)**

$$P(B_7 = 1|Z = 64) = 0.74$$

**(c)**    As shown in Fig. 1, we plot estimated probability every $2 \times 10^4$ samples and we have $1 \times 10^6$ samples in total. And we can tell that our estimate has converged to a good degree of precision (two significant digits).

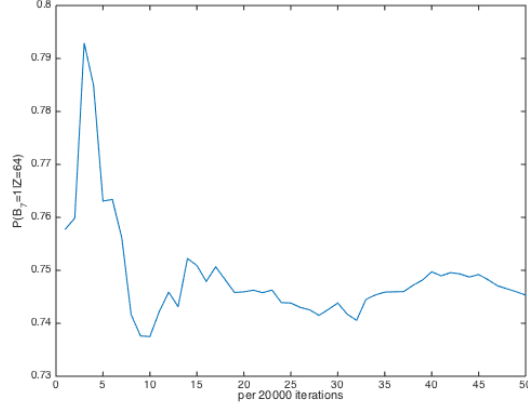## 3.3 Node clustering

CPTs for the polytree is shown in Tab. 1.

Figure 1: $P(B_7 = 1|Z = 64)$ as a function of the number of samples

| $Y_1$ | $Y_2$ | $Y_3$ | $Y$ | $P(Y|X=0)$ | $P(Y|X=1)$ | $P(Z_1=1|Y)$ | $P(Z_2=1|Y)$ |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0.09375 | 0.09375 | 0.9 | 0.1 |
| 1 | 0 | 0 | 2 | 0.28125 | 0.09375 | 0.8 | 0.2 |
| 0 | 1 | 0 | 3 | 0.09375 | 0.03125 | 0.7 | 0.3 |
| 0 | 0 | 1 | 4 | 0.03125 | 0.28125 | 0.6 | 0.4 |
| 1 | 1 | 0 | 5 | 0.28125 | 0.03125 | 0.5 | 0.5 |
| 1 | 0 | 1 | 6 | 0.09375 | 0.28125 | 0.4 | 0.6 |
| 0 | 1 | 1 | 7 | 0.03125 | 0.09375 | 0.3 | 0.7 |
| 1 | 1 | 1 | 8 | 0.09375 | 0.09375 | 0.2 | 0.8 |

Table 1: CPTs for the polytree

## 3.4 Maximum likelihood estimation

**(a)**

$$P(X_{t+1} = x'|X_t = x) = \frac{COUNT_t(x, x')}{COUNT_t(x)} \quad 1 \le t < T$$

**(b)**

$$P(X_t = x|X_{t+1} = x') = \frac{COUNT_t(x, x')}{COUNT_{t+1}(x')} \quad 1 \le t < T$$

**(c)** We first derive joint distribution from G1:

$$
\begin{aligned}
P(X_1 = x_1, X_2 = x_2, ..., X_T = x_T) &= P(X_1 = x_1)\prod_{i=2}^{T} P(X_i = x_i|X_1 = x_1, ..., X_{i-1} = x_{i-1}) \quad (productrule) \\
&= P(X_1 = x_1)\prod_{i=2}^{T} P(X_i = x_i|X_{i-1} = x_{i-1}) \quad (d-separationI) \\
&= \frac{COUNT_1(x_1)}{|data|} \prod_{i=1}^{T-1} \frac{COUNT_i(x_i, x_{i+1})}{COUNT_i(x_i)} \\
&= \frac{\prod_{i=1}^{T-1} COUNT_i(x_i, x_{i+1})}{|data| \prod_{i=2}^{T-1} COUNT_i(x_i)}
\end{aligned}
$$

3

| Token | $P_u(w)$ |
|---|---|
| MILLION | 0.002073 |
| MORE | 0.001709 |
| MR. | 0.001442 |
| MOST | 0.000788 |
| MARKET | 0.000780 |
| MAY | 0.000730 |
| M. | 0.000703 |
| MANY | 0.000697 |
| MADE | 0.000560 |
| MUCH | 0.000515 |
| MAKE | 0.000514 |
| MONTH | 0.000445 |
| MONEY | 0.000437 |
| MONTHS | 0.000406 |
| MY | 0.000400 |
| MONDAY | 0.000382 |
| MAJOR | 0.000371 |
| MILITARY | 0.000352 |
| MEMBERS | 0.000336 |
| MIGHT | 0.000274 |
| MEETING | 0.000266 |
| MUST | 0.000267 |
| ME | 0.000264 |
| MARCH | 0.000260 |
| MAN | 0.000253 |
| MS. | 0.000239 |
| MINISTER | 0.000240 |
| MAKING | 0.000212 |
| MOVE | 0.000210 |
| MILES | 0.000206 |

Table 2: Tokens that start with the letter M and their numerical unigram probabilities

where $|data|$ is the size of the data set. Then we derive joint distribution from G2:

$$
\begin{aligned}
P(X_1 = x_1, X_2 = x_2, ..., X_T = x_T) &= P(X_T = x_T) \prod_{i=1}^{T-1} P(X_i = x_i | X_{i+1} = x_{i+1}, ..., X_T = x_T) \quad (product rule) \\
&= P(X_T = x_T) \prod_{i=1}^{T-1} P(X_i = x_i | X_{i+1} = x_{i+1}) \quad (d-separation I) \\
&= \frac{COUNT_T(x_T)}{|data|} \prod_{i=1}^{T-1} \frac{COUNT_i(x_i, x_{i+1})}{COUNT_{i+1}(x_{i+1})} \\
&= \frac{\prod_{i=1}^{T-1} COUNT_i(x_i, x_{i+1})}{|data| \prod_{i=2}^{T-1} COUNT_i(x_i)}
\end{aligned}
$$

We find the derived joint distribution from G1 and G2 is the same.

## 3.5 Statistical language modeling

**(a)** The results are displayed in Tab. 2.

**(b)** The ten most likely words to follow the word THE, along with their numerical bigram probabilities, are shown in Tab. 3.

| Token | $P_b(w|\textbf{THE})$ |
|---|---|
| ⟨UNK ⟩ | 0.615020 |
| U. | 0.013372 |
| FIRST | 0.011720 |
| COMPANY | 0.011659 |
| NEW | 0.009451 |
| UNITED | 0.008672 |
| GOVERNMENT | 0.006803 |
| NINETEEN | 0.006651 |
| SAME | 0.006287 |
| TWO | 0.006161 |

Table 3: 10 most probable token after **THE**

**(c)**

$$
\begin{aligned}
\mathcal{L}_u &= \log[P_u(\textbf{the})P_u(\textbf{stock})P_u(\textbf{market})...P_u(\textbf{points})P_u(\textbf{last})P_u(\textbf{week})] \\
&= -64.5094403436 \\
\mathcal{L}_b &= \log[P_b(\textbf{the}|\langle\textbf{s}\rangle)P_b(\textbf{stock}|\textbf{the})P_b(\textbf{market}|\textbf{stock})...P_b(\textbf{last}|\textbf{points})P_b(\textbf{week}|\textbf{last})] \\
&= -40.9181321338
\end{aligned}
$$

Bigram model yields the highest log-likelihood.

**(d)**

$$
\begin{aligned}
\mathcal{L}_u &= \log[P_u(\textbf{the})P_u(\textbf{sixteen})P_u(\textbf{officials})...P_u(\textbf{sold})P_u(\textbf{fire})P_u(\textbf{insurance})] \\
&= -44.2919344731 \\
\mathcal{L}_b &= \log[P_b(\textbf{the}|<\textbf{s}>)P_b(\textbf{sixteen}|\textbf{the})P_b(\textbf{officials}|\textbf{sixteen})...P_b(\textbf{fire}|\textbf{sold})P_b(\textbf{insurance}|\textbf{fire})] \\
&= \log(0.0)
\end{aligned}
$$

When the pairs (**sisteen, officials**) and (**sold, fire**) are not observed in the training corpus? This makes the estimated probability to log(0), which is meaningless.

**(e)** Fig. 2 shows the value of this log-likelihood $\mathcal{L}_m$ as a function of the parameter $\lambda \in [0,1]$. And the optimal $\lambda$ is 0.65 with probability -42.9642.
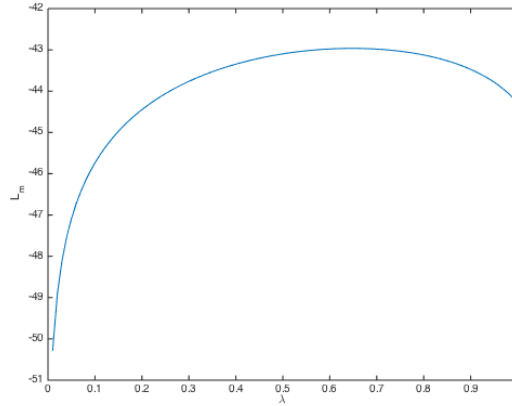


Figure 2: log-likelihood function $\mathcal{L}_m$