# 4.1 Maximum likelihood estimation

## (a)

$$P(Y=y) = \frac{count(Y=y)}{T}$$

$$P(X=x \mid Y=y) = \frac{count(X=x, Y=y)}{\sum_{x'} count(X=x', Y=y)}$$

$$P(Z=z \mid Y=y) = \frac{count(Y=y, Z=z)}{\sum_{z'} count(Y=y, Z=z')}$$

## (b)

$$P(Z=z) = \frac{count(Z=z)}{T}$$

$$P(Y=y \mid Z=z) = \frac{count(Y=y, Z=z)}{\sum_{y'} count(Y=y', Z=z)}$$

$$P(X=x \mid Y=y) = \frac{count(X=x, Y=y)}{\sum_{x'} count(X=x', Y=y)}$$

## (c)

in (a):

$$P(X=x, Y=y, Z=z) = P(Y=y) P(X=x \mid Y=y) P(Z=z \mid Y=y)$$

$$= \frac{count(Y=y)}{T} \frac{count(X=x, Y=y)}{\sum_{x'} count(X=x', Y=y)} \frac{count(Y=y, Z=z)}{\sum_{z'} count(Y=y, Z=z')}$$

$$= \frac{count(X=x, Y=y) \, count(Y=y, Z=z)}{T \, count(Y=y)}$$

in (b):

$$P(X=x, Y=y, Z=z) = P(Z=z)P(Y=y \mid Z=z)P(X=x \mid Y=y)$$

$$= \frac{\cancel{count(Z=z)}}{T} \cdot \frac{count(Y=y, Z=z)}{\cancel{\sum_{y'} count(Y=y', Z=z)}} \cdot \frac{count(X=x, Y=y)}{\sum_{x'} count(X=x', Y=y)}$$

$$= \frac{count(X=x, Y=y) \; count(Y=y, Z=z)}{T \; count(Y=y)}$$

therefore:

they have same joint distribution over $X$, $Y$ and $Z$.

(d)

No. Because

$$P(\vec{F}=\vec{f} \mid \vec{C}=\vec{c}) = \frac{P(\vec{F}=\vec{f}, \vec{C}=\vec{c})}{P(\vec{C}=\vec{c})}$$

$$= \frac{\sum_{(X,Y,Z) \setminus \{\vec{F}, \vec{c}\}} P(X=x, Y=y, Z=z)}{\sum_{(X,Y,Z) \setminus \{\vec{c}\}} P(X=x, Y=y, Z=z)}$$

Because in part(C) we prove that $P(X=x, Y=y, Z=z)$ maybe same. The $P(\vec{F}=\vec{f} \mid \vec{C}=\vec{c})$ maybe no relevant to DAGs.

## 4.2 Markov modeling

(a)

$$P_u(S) = \prod_{i=1}^{L} P_1(z_i)$$

Log-likelihood:

$$L(S) = \sum_{i=1}^{L} \log P_1(z_i)$$

$$= \sum_{\pi} count(\pi) \log P_1(\pi)$$

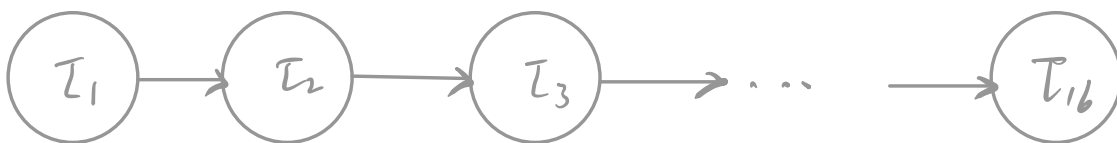Using the answer of HW1, when

$$\frac{count(\pi)}{P_1(\pi)} = Const, \text{ then } L \text{ get largest}$$

number.

| $z$ | a | b | c | d | |
|-----|---|---|---|---|---|
| $P_1(z)$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{4}$ | |

(b)

The BN is as follows:

Because all variables are visible, we can use ML to get the answer.

$$P_1(\tau_1) = \text{count}(\tau_1)$$

$$P_2(\tau' | \tau) = \frac{\text{count}(\tau_i, \tau_{i-1})}{\sum_{\tau_i'} \text{count}(\tau_i', \tau_{i-1})} \quad (i = 2, \cdots, 16)$$

| $P_2'(\tau'|\tau)$ | a | b | c | d |
|---|---|---|---|---|
| a | aa $\frac{1}{2}$ | ab $\frac{1}{4}$ | ac $0$ | ad $\frac{1}{4}$ |
| b | ba $0$ | bb $\frac{3}{4}$ | bc $\frac{1}{4}$ | bd $0$ |
| c | ca $0$ | cb $0$ | cc $\frac{2}{3}$ | cd $\frac{1}{3}$ |
| d | da $\frac{1}{4}$ | db $0$ | dc $\frac{1}{4}$ | dd $\frac{1}{2}$ |

(c)
The possible form is contained in the (b) table, Obviously the $\tau_2 \cdot \tau_3$ underlined string are 0 in the table

We may find that in $S, T_1, T_2, T_3$, the number of every letter is same. Therefore, is the number of letters is $n$, then obviously:

$$P_u(T) = (\frac{1}{n})^L$$

Therefore, when $n$ increases, then $P_u(T)$ decreases. In conclusion,

$$P_u(S) < P_u(T_1)$$
$$P_u(S) = P_u(T_2)$$
$$P_u(S) < P_u(T_3)$$

$P_B(T)$ is relevant to the kinds of two-letter string. If there are more kinds of it, the $P_B(T)$ will be smaller.
In conclusion.

$$P_B(T_1) > P_B(S)$$
$$P_B(T_2) > P_B(S)$$
$$P_B(T_3) > P_B(T_2)$$

Because the factor of $P_u(S), P_u(T_2)$ is $\frac{1}{4}$, and factor of $P_B(S), P_u(T_2)$ is bigger than $\frac{1}{4}$. Therefore,

$P_u(S) < P_B(S)$

$P_u(T_2) < P_B(T_2)$

Similarly, the factor of $P_u(T_1)$ and $P_u(T_3)$ is $1$, and the factor of $P_B(T_1)$ and $P_B(T_3)$ is $\frac{1}{2}$. Therefore,

$P_u(T_1) < P_B(T_1)$

$P_u(T_3) < P_B(T_3)$