

General Mixture

written by Khashayar

February 2018

1 Problem formulation

We have N number of binary mother sequences $\mathcal{Y} = \{y_i\}_{i=1}^N$ with length L , each of them have a corresponding probability f_i , where $\sum_{i=1}^N f_i = 1$. In each time step $1 \leq t \leq m$, one of the mother sequences such as y_{r_t} is picked with respect to the probabilities $\mathcal{F} = \{f_i\}$, and is passed through a symmetric noisy channel with flip probability f . In the other end, we observe the noisy samples $\{z_i\}_{i=1}^m$. We know the channel's flip rate f and the sequences' length L , and we want to estimate N , the mother sequences themselves, and their frequencies $\{f_i\}_{i=1}^N$. We will introduce an estimator and bound it's norm-2 error using notions from high dimensional statistics.

2 The Estimator

We denote the space of all binary sequences with length L by \mathcal{A}_L , which has size $n = 2^L$. By \mathcal{A}_L^m and \mathcal{A}_L^s we refer to the space of mother sequences and sample sequences respectively. Let $p = \{p_i\}_{i=1}^n$ be the probability distribution induced by the set of mother sequences \mathcal{Y} on \mathcal{A}_L^s . Furthermore, for sequences s_1, s_2 , let $d(s_1, s_2)$ be their hamming distance, and $p(s_1, s_2) = f^{d(s_1, s_2)}(1 - f)^{L - d(s_1, s_2)}$.

2.1 Preliminaries

We denote the empirical distribution on \mathcal{A}_L^s by $q = \{q_i\}_{i=1}^n$, which is obtained by counting the number of observed sequences for each type in \mathcal{A}_L^s , and normalizing them by their sum m . Look at q as a random distribution, thus each q_i is a random variable. Obviously, we see that $\mathbb{E}[q] = p$. In fact, $\{q_i\}_{i=1}^n$ is the normalized version of a multinomial distribution with probabilities $\{p_i\}_{i=1}^n$. Hence

$$Cov[q_i, q_j] = -\frac{p_i p_j}{m}, \quad Var[p_i] = \frac{p_i(1 - p_i)}{m}, \quad \forall 1 \leq i < j \leq n. \quad (1)$$

Let $\epsilon = \{\epsilon_i\}_{i=1}^n$ be $\epsilon = p - q$, for which $\mathbb{E}[\epsilon] = 0$. Equation (1) gives

$$\mathbb{E}[\epsilon_i \epsilon_j] = -\frac{p_i p_j}{m}, \quad \mathbb{E}[\epsilon_i^2] = \frac{p_i(1-p_i)}{m}, \quad \forall 1 \leq i < j \leq n. \quad (2)$$

Now, Let B be a $n \times n$ matrix which has its rows and columns corresponded to sequences in \mathcal{A}_L , such that for every sequences $s_1, s_2 \in \mathcal{A}_L$,

$$B(s_1, s_2) = p(s_1, s_2),$$

where $B(s_1, s_2)$ is the matrix element corresponding to row s_1 and column s_2 . Let $\mathcal{Q}(\mathcal{A}_L)$ be the space of probability distributions on \mathcal{A}_L . For desired $v \in \mathcal{Q}(\mathcal{A}_L)$ and $s \in \mathcal{A}_L$, the element in vector v which corresponds to sequence s is denoted by v_s . It can be easily checked that B is an invertible probability matrix, hence it defines an onto and one to one linear transformation from $\mathcal{Q}(\mathcal{A}_L^m)$ to $\mathcal{Q}(\mathcal{A}_L^s)$. Now We apply the inverse map B^{-1} to the empirical distribution q in $\mathcal{Q}(\mathcal{A}_L^s)$, to obtain the empirical distribution $Y = \{Y_s\}_{s \in \mathcal{A}_L^m}$ in $\mathcal{Q}(\mathcal{A}_L^m)$. Let \mathcal{F} be a vector in $\mathcal{Q}(\mathcal{A}_L^m)$ such that for each $s \in \mathcal{A}_L$

$$\mathcal{F}_s = \begin{cases} f_i & s = y_i \in \mathcal{Y} \\ 0 & o.w. \end{cases}$$

By the above definition, our goal is to estimate \mathcal{F} . Note that q is a noisy unbiased observation of p , and it can be readily seen that $p = B(\mathcal{F})$. Noting the fact that B^{-1} is linear, we deduce that

$$\mathbb{E}[Y] = \mathbb{E}[B^{-1}(q)] = B^{-1}(\mathbb{E}[q]) = B^{-1}(p) = \mathcal{F}. \quad (3)$$

Hence, if we define $\alpha \in \mathcal{Q}(\mathcal{A}_L^m)$ such that for each $s \in \mathcal{A}_L^m$,

$$\alpha_s = Y_s - \mathcal{F}_s,$$

Then according to equation (3), Y is an unbiased noisy observation of \mathcal{F} with noise α . Thus $\mathbb{E}[\alpha] = 0$. In addition, by definition of ϵ , observe that

$$\alpha = B^{-1}(\epsilon). \quad (4)$$

2.2 Calculation of Noise Variance

Now, we are interested in calculating the errors $\mathbb{E}[\alpha_s]$ for each $s \in \mathcal{A}_L^m$.

It can be seen that B is the result of tensoring the matrix $\begin{bmatrix} 1-f & f \\ f & 1-f \end{bmatrix}$ to itself L times. Hence,

$$B^{-1} = \frac{1}{(1-2f)^L} \underbrace{\begin{bmatrix} 1-f & -f \\ -f & 1-f \end{bmatrix} \otimes \cdots \otimes \begin{bmatrix} 1-f & -f \\ -f & 1-f \end{bmatrix}}_{n \text{ times}}. \quad (5)$$

By equation (5) we obtain that

$$B^{-1}(s_1, s_2) = D_f \cdot (-1)^{d(s_1, s_2)} p(s_1, s_2), \quad (6)$$

where $D_f = \frac{1}{(1-2f)^L}$.

Remark that we can index the elements of p, q, ϵ by sequences in \mathcal{A}_L as well.

According to equations (4) and (6), for each $s^* \in \mathcal{A}_L^m$,

$$\alpha_{s^*} = D_f \sum_{s \in \mathcal{A}_L^s} p(s^*, s) \epsilon_s (-1)^{d(s^*, s)}. \quad (7)$$

Hence,

$$\begin{aligned} \mathbb{E}[\alpha_{s^*}^2] &= \mathbb{E}[D_f^2 (\sum_{s \in \mathcal{A}_L^s} p(s^*, s) \epsilon_s (-1)^{d(s^*, s)})^2] \\ &= D_f^2 \sum_{s \in \mathcal{A}_L^s} \mathbb{E}[p(s^*, s)^2 \epsilon_s^2] + 2D_f^2 \sum_{s_1 \neq s_2 \in \mathcal{A}_L^s} \mathbb{E}[p(s^*, s_1) p(s^*, s_2) \epsilon_{s_1} \epsilon_{s_2} (-1)^{d(s_1, s_2)}] \\ &= \frac{D_f^2}{m} \sum_{s \in \mathcal{A}_L^s} p(s^*, s)^2 p_s (1 - p_s) - \frac{2D_f^2}{m} \sum_{s_1 \neq s_2 \in \mathcal{A}_L^s} p(s^*, s_1) p(s^*, s_2) p_{s_1} p_{s_2} (-1)^{d(s_1, s_2)} \\ &= R_f [\sum_{s \in \mathcal{A}_L^s} p(s^*, s)^2 p_s - (\sum_{s \in \mathcal{A}_L^s} p(s^*, s) p_s (-1)^{d(s^*, s)})^2], \end{aligned} \quad (8)$$

where

$$R_f = \frac{1}{(1-2f)^{2L} m}. \quad (9)$$

Note that $p = p_{s \in \mathcal{A}_L^s}$ is the distribution induced by the mother sequences on \mathcal{A}_L^s .

$$p_s = \sum_{i=1}^N f_i p(y_i, s). \quad (10)$$

Thus (8) gives

$$\mathbb{E}[\alpha_{s^*}^2] \quad (11)$$

$$= R_f \left[\sum_{s \in \mathcal{A}_L^s} p(s^*, s)^2 \left(\sum_{i=1}^N f_i p(y_i, s) \right) - \left[\sum_{s \in \mathcal{A}_L^s} p(s^*, s) \left(\sum_{i=1}^N f_i p(y_i, s) \right) (-1)^{d(s^*, s)} \right]^2 \right] \quad (12)$$

$$= R_f \left[\sum_{i=1}^N f_i \sum_{s \in \mathcal{A}_L^s} p(s^*, s)^2 p(y_i, s) - \left[\sum_{s \in \mathcal{A}_L^s} p(s^*, s) \left(\sum_{i=1}^N f_i p(y_i, s) \right) (-1)^{d(s^*, s)} \right]^2 \right] \quad (13)$$

$$= R_f \left[\sum_{i=1}^N f_i \sum_{s \in \mathcal{A}_L^s} p(s^*, s)^2 p(y_i, s) - \left[\sum_{i=1}^N f_i \sum_{s \in \mathcal{A}_L^s} p(s^*, s) p(y_i, s) (-1)^{d(s^*, s)} \right]^2 \right] \quad (14)$$

$$= R_f \left[\left(\sum_{i=1}^N f_i (f^3 + (1-f)^3)^{L-d(s^*, y_i)} (f(1-f))^{d(s^*, y_i)} \right) - \psi(s^*) \right], \quad (15)$$

$$(16)$$

$$\text{Where } \psi(s^*) = \begin{cases} f_i^2 (1-2f)^{2L} & s^* = y_i \in \mathcal{Y} \\ 0 & o.w. \end{cases}$$

We observe that if f is relatively small, then the major mass of noise is on the mother sequences themselves.

2.3 Estimation

For a subset κ of sequences in \mathcal{A}_L^m , define $\hat{\mathcal{F}}_\kappa$ such that for each sequence $s \in \mathcal{A}_L^m$,

$$\hat{\mathcal{F}}_{\kappa s} = \begin{cases} Y_s & s \in \kappa \\ 0 & o.w. \end{cases}$$

Define

$$\hat{\mathcal{F}} = \hat{\mathcal{F}}_{\hat{\kappa}} = \arg \min_{\kappa} \{ \| \hat{\mathcal{F}}_{\kappa} - Y \|^2 + \text{pen}(\kappa) \}. \quad (17)$$

At the moment, we define

$$\text{pen}(\kappa) = R_f [(f^3 + (1-f)^3)^L + (|\kappa| - 1)f(1-f)(f^3 + (1-f)^3)^{L-1}]. \quad (18)$$

However, we will show later that we can use a more complicated penalty term which gives a better accuracy, with the cost of a higher computational cost. Also, note that in equation (18) the first term is a constant and does not depend on κ , thus has no effect on our model selection, but for sake of simplicity of the proof, we don't omit it.

Remark that if we sort the elements of vector Y , it is sufficient to evaluate the minimization problem in (17) only on the top r sequences with the largest values in Y . Therefore, the computational cost of our estimator is linear.

2.4 Error Bound

Note that by definition, $\mathcal{F}_{\mathcal{Y}}$ is a vector equal to \mathbf{Y} with respect to the mother sequences, and zero otherwise. Now according to equation (17), we can write

$$\|\hat{\mathcal{F}} - \mathbf{Y}\|^2 + \text{pen}(\hat{\kappa}) \leq \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathbf{Y}\|^2 + \text{pen}(\mathcal{Y}),$$

which gives

$$\begin{aligned} \|\hat{\mathcal{F}} - \mathcal{F}\|^2 + \|\mathcal{F} - \mathbf{Y}\|^2 + 2\langle \hat{\mathcal{F}} - \mathcal{F}, \mathcal{F} - \mathbf{Y} \rangle + \text{pen}(\hat{\kappa}) &\leq \\ \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2 + \|\mathcal{F} - \mathbf{Y}\|^2 + 2\langle \hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}, \mathcal{F} - \mathbf{Y} \rangle + \text{pen}(\mathcal{Y}). \end{aligned} \quad (19)$$

But note that for $s \notin \mathcal{Y}$,

$$\mathbb{E}[(\hat{\mathcal{F}}_{\mathcal{Y}_s} - \mathcal{F}_s)(\mathcal{F}_s - \mathbf{Y}_s)] = \mathbb{E}[-\mathcal{F}_s(\mathcal{F}_s - \mathbf{Y}_s)] = 0,$$

and remark that $\mathcal{F} - \mathbf{Y} = \epsilon$, hence

$$\langle \hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}, \mathcal{F} - \mathbf{Y} \rangle = \langle \alpha_{\mathcal{Y}}, -\alpha \rangle = -\|\alpha_{\mathcal{Y}}\|^2 \leq 0,$$

where $\epsilon_{\mathcal{Y}}$ is the projection of α on the indices with respect to the mother sequences. In the same manner, observe that

$$\langle \hat{\mathcal{F}} - \mathcal{F}, \mathcal{F} - \mathbf{Y} \rangle = \langle \alpha_{\hat{\kappa}}, -\alpha \rangle = -\|\alpha_{\hat{\kappa}}\|^2. \quad (20)$$

Thus, equation (19) gives

$$\|\hat{\mathcal{F}} - \mathcal{F}\|^2 \leq \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2 + \text{pen}(\mathcal{Y}) + \|\alpha_{\hat{\kappa}}\|^2 - \text{pen}(\hat{\kappa}),$$

which gives

$$\mathbb{E} \|\hat{\mathcal{F}} - \mathcal{F}\|^2 \leq \mathbb{E} \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2 + \text{pen}(\mathcal{Y}) + \mathbb{E} [\|\alpha_{\hat{\kappa}}\|^2 - \text{pen}(\hat{\kappa})]. \quad (21)$$

But due to equation (15), for a fixed κ we have

$$\begin{aligned}
\mathbb{E} \|\alpha_\kappa\|^2 &= \sum_{s^* \in \kappa} R_f \left[\left(\sum_{i=1}^N f_i (f^3 + (1-f)^3)^{L-d(s^*, y_i)} (f(1-f))^{d(s^*, y_i)} \right) - \psi(s^*) \right] \\
&\leq R_f \sum_{s^* \in \kappa} \left(\sum_{i=1}^N f_i (f^3 + (1-f)^3)^{L-d(s^*, y_i)} (f(1-f))^{d(s^*, y_i)} \right) \\
&= R_f \sum_{s^* \in \kappa \cap \mathcal{Y}} \left(\sum_{i=1}^N f_i (f^3 + (1-f)^3)^{L-d(s^*, y_i)} (f(1-f))^{d(s^*, y_i)} \right) \\
&\quad + R_f \sum_{\substack{s^* \in \kappa \\ s^* \notin \mathcal{Y}}} \left(\sum_{i=1}^N f_i (f^3 + (1-f)^3)^{L-d(s^*, y_i)} (f(1-f))^{d(s^*, y_i)} \right) \\
&\leq R_f \left(\sum_{s^* = y_j \in \kappa \cap \mathcal{Y}} f_j (f^3 + (1-f)^3)^L + \sum_{i \neq j} f_i (f^3 + (1-f)^3)^{L-1} f(1-f) \right) \\
&\quad + R_f \sum_{\substack{s^* \in \kappa \\ s^* \notin \mathcal{Y}}} \sum_{i=1}^N f_i (f^3 + (1-f)^3)^{L-1} f(1-f) \\
&\leq R_f [(f^3 + (1-f)^3)^L + (|\kappa| - 1)(f^3 + (1-f)^3)^{L-1} f(1-f)] = \text{pen}(\kappa). \tag{22}
\end{aligned}$$

On the other hand,

$$\begin{aligned}
\mathbb{E} \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2 &= \mathbb{E} \|\alpha_{\mathcal{Y}}\|^2 \tag{23} \\
&= \sum_{y_j \in \mathcal{Y}} R_f \left[\left(\sum_{i=1}^N f_i (f^3 + (1-f)^3)^{L-d(y_j, y_i)} (f(1-f))^{d(y_j, y_i)} \right) - \psi(y_j) \right] \\
&\geq R_f \sum_{y_j \in \mathcal{Y}} [(f_j (f^3 + (1-f)^3)^{L-d(y_j, y_j)} (f(1-f))^{d(y_j, y_j)}) - \psi(y_j)] \\
&= R_f \left[\left(\sum_{i=1}^N f_i \right) (f^3 + (1-f)^3)^L - \left(\sum_{i=1}^N f_i^2 \right) (1-2f)^{2L} \right]. \tag{24}
\end{aligned}$$

But it can easily be checked that

$$f^3 + (1-f)^3 \geq (1-2f)^2,$$

thus, we obtain that

$$\mathbb{E} \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2 \geq R_f [(f^3 + (1-f)^3)^L - \left(\sum_{i=1}^N f_i^2 \right) (f^3 + (1-f)^3)^L]. \tag{25}$$

Now suppose that we have a constraint for $0 \leq \gamma \leq 1$ as

$$\min_i f_i \geq \frac{\gamma}{N},$$

which is essential so that the problem becomes identifiable. Then, we get

$$\sum_i f_i^2 \geq \frac{(N-1)\gamma^2}{N^2} + (1 - \frac{(N-1)\gamma}{N})^2 = 1 - \frac{N-1}{N}\gamma(2-\gamma). \quad (26)$$

inequalities (25) , (26) give

$$\begin{aligned} \mathbb{E} \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2 &\geq \frac{N-1}{N}\gamma(2-\gamma)R_f(f^3 + (1-f)^3)^L \\ &= \frac{N-1}{N}\gamma(2-\gamma)\frac{1}{1 + \frac{f(1-f)}{f^3+(1-f)^3}(N-1)}\text{pen}(\mathcal{Y}). \end{aligned} \quad (27)$$

Define the ratio

$$K_{f,\gamma,N} = \frac{N}{N-1} \frac{\left(1 + \frac{f(1-f)}{f^3+(1-f)^3}(N-1)\right)}{\gamma(2-\gamma)}.$$

Then, According to inequality (27),

$$\text{pen}(\mathcal{Y}) \leq K_{f,\gamma,N} \mathbb{E} \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2. \quad (28)$$

Now observe that if $\hat{\kappa}$ was fixed, then inequality (22) gives

$$\mathbb{E} [\|\alpha_{\hat{\kappa}}\|^2 - \text{pen}(\hat{\kappa})] \leq 0.$$

Consequently, due to inequality (21),

$$\mathbb{E} \|\hat{\mathcal{F}} - \mathcal{F}\|^2 \leq \mathbb{E} \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2 + \text{pen}(\mathcal{Y}), \quad (29)$$

which combined by (28) reveals

$$\mathbb{E} \|\hat{\mathcal{F}} - \mathcal{F}\|^2 \leq (1 + K_{f,\gamma,N}) \mathbb{E} \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2. \quad (30)$$

Moreover, we can find an explicit bound for the expected norm-2 error $\mathbb{E} \|\hat{\mathcal{F}} - \mathcal{F}\|^2$ by driving a upper bound for $\mathbb{E} \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2$. According to the inequality (22), we have

$$\text{pen}(\mathcal{Y}) \geq \|\alpha_{\mathcal{Y}}\|^2 = \mathbb{E} \|\hat{\mathcal{F}}_{\mathcal{Y}} - \mathcal{F}\|^2. \quad (31)$$

Inequalities (31) , (29) reveal

$$\begin{aligned} \mathbb{E} \|\hat{\mathcal{F}} - \mathcal{F}\|^2 &\leq \\ 2\text{pen}(\mathcal{Y}) &= 2R_f[(f^3 + (1-f)^3)^L + (N-1)f(1-f)(f^3 + (1-f)^3)^{L-1}]. \end{aligned} \quad (32)$$