

DEADLINE: 13.06.2020, 23:59 (Moodle)

Points: 20 (obligatory part) + 4 (bonus part)

Motif finding in DNA sequences

1 Introduction

TO HAND OUT: A report in a .pdf file and working files in Python (ver 3).

The report should be written, so that someone not attending the lecture can understand it. Thus, include formulation of the problem etc.

REMARK: The description of what should be included in the project is provided below. Your project should contain all these elements, but it does not restrict itself only to them.

Project may be performed in teams, at most 2 persons are allowed in one team.

Remark: The project is related to the EM algorithm, it is

- on the one hand: an extension of *flipping a coin* example presented in the lecture,
- on the other hand: it is a simplified model related to so-called *motif finding* in DNA sequences.

You can get up to **20 pts** for an obligatory part and up to **4 pts** for a “bonus” part.

2 Applications of the EM algorithm to a simplified version of motif finding in DNA sequences.

2.1 Description of the problem

Remark: The problem itself can be stated in a (significantly) shorter version. I present it together with some introduction.

Assume we have a sequence $\mathbf{x}_1 = (x_{11}, x_{12}, \dots, x_{1w})$ of length w , where each $x_{1i} \in \Sigma = \{A, C, G, T\}$. For a convenience, we will identify¹ $\{A, C, G, T\} \equiv \{1, 2, 3, 4\}$. Denote

¹Writing e.g., “a letter 2” we mean a letter C, etc.

$d = |\Sigma|$, i.e., in our case $d = 4$. We are given distributions $\boldsymbol{\theta}_j = (\theta_{1,j}, \dots, \theta_{d,j})^T$, $j = 1, \dots, w$. The sequence \mathbf{x}_1 is a realization of a random variable $X = (X_1, \dots, X_w)$, where

$$P(X_j = a) = \theta_{a,j}, a = 1, \dots, d \quad (\equiv a \in \{A, C, G, T\})$$

In other words $\theta_{a,j}$ is the probability that the letter a is present on j -th position. In biological sciences In biological sciences the matrix

$$\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_w)$$

is often called “a position weight matrix of a motif”. Thus, having the matrix $\boldsymbol{\theta}$, the probability of observing $\mathbf{x}_1 = (x_{11}, x_{12}, \dots, x_{1w})$ is equal to

$$P(X = \mathbf{x}_1; \boldsymbol{\theta}) \equiv P(\mathbf{x}_1; \boldsymbol{\theta}) = \prod_{i=1}^w \theta_{x_{1i}, i}$$

Przykład 2.1.1 Assume that $w = 3$ and $\boldsymbol{\theta}_1 = (3/8, 1/8, 2/8, 2/8)^T$, $\boldsymbol{\theta}_2 = (1/10, 2/10, 3/10, 4/10)^T$, $\boldsymbol{\theta}_3 = (1/7, 2/7, 1/7, 3/7)^T$. The probability of observing $\mathbf{x}_1 = (T, G, A) \equiv (4, 3, 1)$ is equal to

$$P(\mathbf{x}_1; \boldsymbol{\theta}) = P((T, G, A); \boldsymbol{\theta}) = P((4, 3, 1); \boldsymbol{\theta}) = \theta_{1,4} \cdot \theta_{2,3} \cdot \theta_{3,1} = \frac{2}{8} \cdot \frac{3}{10} \cdot \frac{1}{7}.$$

Modification 1. Assume now that we have k observations (aka rows) $\mathbf{x}_1, \dots, \mathbf{x}_k$: (each \mathbf{x}_i was generated in a way described above)

$$\begin{aligned} \mathbf{x}_1 &= (x_{11}, x_{12}, \dots, x_{1w}) \\ \mathbf{x}_2 &= (x_{21}, x_{22}, \dots, x_{2w}) \\ &\vdots \\ \mathbf{x}_k &= (x_{k1}, x_{k2}, \dots, x_{kw}) \end{aligned}$$

The probability of observing $\mathbf{x}_1, \dots, \mathbf{x}_k$ may be written as:

$$L(\boldsymbol{\theta}) = P(\mathbf{x}_1, \dots, \mathbf{x}_k; \boldsymbol{\theta}) = \prod_{j=1}^k \prod_{i=1}^w \theta_{x_{ji}, i}. \quad (1)$$

Modification 2. Assume now that we have observed $\mathbf{x}_1, \dots, \mathbf{x}_k$, but we **do not know** $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_w$ (we still know w and α). How to estimate these parameters? MLE is

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}).$$

It is easy to check that then

$$\hat{\theta}_{ij} = \frac{k_{ij}}{k},$$

where k_{ij} is a number of letters j on positions i in all k sequences.

Modification 3. Now we complicate things. Fix $\alpha \in (0, 1)$ (in this – obligatory – part we assume that α is **known**). Assume that $\mathbf{x}_1, \dots, \mathbf{x}_k$ were obtained in the following way:

- For each $i = 1, \dots, k$ we “flip independently a coin“, i.e., we simulate $Z_i \in \{0, 1\}$:

$$P(Z_i = 1) = \alpha, \quad P(Z_i = 0) = 1 - \alpha$$

- If $Z_i = 1$ then we generate \mathbf{x}_i as previously;
- If $Z_i = 0$ then we generate $\mathbf{x}_i = (x_{i1}, \dots, x_{iw})$, where x_{ij} are i.i.d. with a distribution².

$$\boldsymbol{\theta}^b = (\theta_1^b, \theta_2^b, \theta_3^b, \theta_4^b) \equiv (\theta_A^b, \theta_C^b, \theta_G^b, \theta_T^b).$$

(i.e., the consecutive letters are generated independently – it is also independent from the position, e.g., the probability of obtaining (G, G, C) is $\theta_G^b \cdot \theta_G^b \cdot \theta_C^b$). In other words

$$P(\mathbf{x}_i; \boldsymbol{\theta}) = \prod_{j=1}^w \theta_{x_{ij}}^b.$$

If we knew (z_1, \dots, z_k) , i.e., which of sequences $\mathbf{x}_1, \dots, \mathbf{x}_k$ it from the distribution $\boldsymbol{\theta}$ and which is from the distribution $\boldsymbol{\theta}^b$, then we could write the likelihood function as

$$L(\boldsymbol{\theta}, \boldsymbol{\theta}^b) = \prod_{i=1}^k [z_i P(\mathbf{x}_i; \boldsymbol{\theta}) + (1 - z_i) P(\mathbf{x}_i; \boldsymbol{\theta}^b)]$$

and we could directly find $\arg \max_{\boldsymbol{\theta}, \boldsymbol{\theta}^b} L(\boldsymbol{\theta}, \boldsymbol{\theta}^b)$. It is easy to check that then: for sequences

for which $z_i = 1$ the parameter $\boldsymbol{\theta}$ is estimated computing frequencies of letters on each position, whereas for sequences for which $z_i = 0$ the parameter $\boldsymbol{\theta}^b$ is estimated computing frequencies of letters (independently of position).

Knowing the distribution $P(Z_i = z_i) = \alpha^{z_i} (1 - \alpha)^{1 - z_i}$ we may write the joint distribution of $\mathbf{x}_1, \dots, \mathbf{x}_k$ and $\mathbf{z} = (z_1, \dots, z_k)$

$$P(\mathbf{x}_1, \dots, \mathbf{x}_k, \mathbf{z}; \boldsymbol{\theta}, \boldsymbol{\theta}^b) = \prod_{i=1}^k P(Z_i = z_i) [z_i P(\mathbf{x}_i; \boldsymbol{\theta}) + (1 - z_i) P(\mathbf{x}_i; \boldsymbol{\theta}^b)]$$

Modification 4 (the actual project problem). Everything as previously, we have observations $\mathbf{x}_1, \dots, \mathbf{x}_k$, **but** we do not know $\mathbf{z} = (z_1, \dots, z_k)$.

Task: estimate $\boldsymbol{\theta}$ and $\boldsymbol{\theta}^b$. Denote all the parameters as

$$\boldsymbol{\Theta} = (\boldsymbol{\theta}, \boldsymbol{\theta}^b).$$

Our likelihood is thus the marginal distribution:

$$\begin{aligned} L(\boldsymbol{\Theta}) = P(\mathbf{x}_1, \dots, \mathbf{x}_k; \boldsymbol{\Theta}) &= \sum_{\mathbf{z}=(z_1, \dots, z_k) \in \{0,1\}^k} P(\mathbf{x}_1, \dots, \mathbf{x}_k, \mathbf{z}; \boldsymbol{\Theta}) \\ &= \prod_{i=1}^k (\alpha P(\mathbf{x}_i; \boldsymbol{\theta}) + (1 - \alpha) P(\mathbf{x}_i; \boldsymbol{\theta}^b)). \end{aligned}$$

² $\boldsymbol{\theta}^b$ is called *the background distribution*

More specifically, we want to find/estimate

$$\hat{\Theta} = \arg \max_{\Theta} L(\Theta) = \arg \max_{\Theta} \mathcal{L}(\Theta), \quad (2)$$

where $\mathcal{L}(\Theta) = \log L(\Theta)$.

2.2 Task – Python’s scripts

2.2.1 Generating data.

To generate data we need the following parameters:

- w = length of a motif
- Position weight matrix of the motif $\theta = (\theta_1, \dots, \theta_w)$, where each θ_j is a column vector of size 4 (we consider only letters $\{A, C, G, T\} \equiv \{1, 2, 3, 4\}$)
- The background distribution $\theta^b = (\theta_1^b, \theta_2^b, \theta_3^b, \theta_4^b) \equiv (\theta_A^b, \theta_C^b, \theta_G^b, \theta_T^b)$
- α – the probability that a row comes from ”motif“
- k – number of rows to simulate

Parameters are stored in a file with extension `.json`. The attached (in Moodle) file `json_save_params.py` saves sample parameters of the file `params_set1.json`

You should **write the script for generating the data** using read parameters. The format is as follows:

```
python3 motif_IndexNr_generate.py --params param_file.json
--output generated_data.json
```

The file reads in the parameters (e.g., from a sample `params_set1.json`), generates the data to the file `generated_data.json`. More precisely – we save data + information on used α .

Detailed format – you should see (and follow it) how it is done in `motif_123456_generate.py` (I suggest simply to update this file – reading `.json` and saving to `.json` is prepared therein)

2.2.2 Estimating parameters $\Theta = (\theta, \theta^b)$ using the EM algorithm.

This is actually the **main (obligatory) part of the project**. We read in data in format produced by `motif_IndexNr_generate.py`, assume they are in `generated_data.json`. In other words we have $\mathbf{x}_1, \dots, \mathbf{x}_k$ and – knowing α – you should estimate $\Theta = (\theta, \theta^b)$ – (2).

You need to **write the script for estimating the parameters** using the EM algorithm. The format:

```
python3 motif_IndexNr_estimate.py --input generated_data.json
--output estimated_params.json
```

Thus, the estimated parameters also should be saved in a `.json` file. Similarly, in Moodle there is a "skeleton" of the above file:

`motif_123456_generate.py`, which reads in data and saves sample (random – need to replace it with your estimations) θ, θ^b in a `.json` file.

How will your programs will be tested? For some w, k, α and θ, θ^b rows $\mathbf{x}_1, \dots, \mathbf{x}_k$ (I can, but need not to, use your `motif_IndexNr_generate.py`) are generated, then the "quality" of your estimates of θ, θ^b will be checked. Parameters w and k will surely be larger than in sample files. More details on evaluation of your estimates are provided in Section ??.

2.3 Bonus task (4 pts) – a parameter α is unknown

This time we are given observations $\mathbf{x}_1, \dots, \mathbf{x}_k$, we thus know w , but we do not know the parameter α (and of course we still do not know θ, θ^b). In other words we are to estimate

$$\theta' = (\theta, \theta^0, \alpha).$$

Technically: a file `motif_IndexNr_estimate.py` takes extra parameter `--estimate-alpha`, which can be either 'no' (default) or 'yes'. Thus, calling

```
python3 motif_IndexNr_estimate.py --input generated_data.json
--output estimated_params.json --estimate-alpha yes
```

should read in `generated_data.json` and ignore the information on α (which is available in the file). Next, all the estimated parameters should be saved in `estimated_params.json`. Note that this way the format of files remains unchanged.

2.4 More details on evaluation of your estimates

For two discrete probability distributions $\mathbf{p} = (p_1, \dots, p_n)^T, \mathbf{q} = (q_1, \dots, q_n)^T$ on $\{1, \dots, n\}$ we define the total variation distance as

$$d_{\text{tv}}(\mathbf{p}, \mathbf{q}) = \frac{1}{2} \sum_{i=1}^n |p_i - q_i|$$

Note that we have $0 \leq d_{\text{tv}}(\mathbf{p}, \mathbf{q}) \leq 1$.

Denote the "original" parameters (the ones used for generating data) as

$$\theta^{b, \text{orig}} = (\theta_1^{b, \text{orig}}, \theta_2^{b, \text{orig}}, \theta_3^{b, \text{orig}}, \theta_4^{b, \text{orig}})$$

$$\theta^{\text{orig}} = (\theta_1^{\text{orig}}, \dots, \theta_w^{\text{orig}})$$

(each $\theta_i^{\text{orig}}, i = 1, \dots, w$ is a column vector of size 4).

Similarly, the estimated parameters are denoted as

$$\theta^{b, \text{estym}} = (\theta_1^{b, \text{estym}}, \theta_2^{b, \text{estym}}, \theta_3^{b, \text{estym}}, \theta_4^{b, \text{estym}})$$

$$\theta^{\text{estym}} = (\theta_1^{\text{estym}}, \dots, \theta_w^{\text{estym}})$$

As a *final performance measure* the following average of total variation distances of all the distributions which appear here will be used, i.e.,

$$d_{\text{tv}} = \frac{1}{w+1} \left[d_{\text{tv}}(\theta^{b, \text{orig}}, \theta^{b, \text{estym}}) + \sum_{i=1}^w d_{\text{tv}}(\theta_i^{\text{orig}}, \theta_i^{\text{estym}}) \right]$$