# PAC learning of Probabilistic Automaton based on the Method of Moments (Supplementary Material)

**Hadrien Glaude**                                                                    HADRIEN.GLAUDE@INRIA.FR

Univ. Lille, CRIStAL, UMR 9189, SequeL Team, Villeneuve d'Ascq, 59650, France

**Olivier Pietquin**[1]                                                          OLIVIER.PIETQUIN@UNIV-LILLE1.FR

Institut Universitaire de France (IUF), Univ. Lille, CRIStAL, UMR 9189, SequeL Team, Villeneuve d'Ascq, 59650, France

## 1. Finite sample analysis of CH-PRFA

In this paper, we prove the following Theorem.

**Theorem 1.** *Let $p$ be a distribution realized by a minimal PRFA of size $d$, $\mathcal{B} = (\mathcal{P}, \mathcal{S})$ be a complete and residual basis, we denote by $\sigma_d$ the $d$-th largest singular values of $(\mathbf{p}_u(v))_{u \in \mathcal{R}}$. Let $\mathcal{D}$ be a training set of words generated by $p$, we denote by $n$ the number of time the least occurring prefix of $\mathcal{P}$ appears in $\mathcal{D}$ ($n = \min_{u \in \mathcal{P}} |\{\exists v \in \Sigma^\star | uv \in \mathcal{D}\}|$). For all $0 < \delta < 1$, there exists a constant $K$ such that, for all $t > 0$, $\epsilon > 0$, with probability $1 - \delta$, if*

$$n \geq K \frac{t^4 d^4 |\Sigma|}{\epsilon^2 \sigma_d^{10}} \log \left( \frac{|\mathcal{P}|}{\delta} \right),$$

*CH-PRFA returns a PFA realizing a proper distribution $\hat{p}$ such that*

$$\sum_{u \in \Sigma^{\leq t}} |\hat{p}(u) - p(u)| \leq \epsilon.$$

### 1.1. Notations

In order to make the proof easier to read, we first define few notations. Let $Z$ be the matrix built by stacking the $\mathbf{p}_u^\top$, where $u \in \mathcal{R}$ as follows,

$$Z = (\mathbf{p}_u)_{u \in \mathcal{R}}^\top.$$

Similarly, we define the following matrices :

$$\hat{Z} = (\mathbf{p}_u)_{u \in \hat{\mathcal{R}}}^\top,$$
$$Z_o = (\dot{o}\mathbf{p}_u)_{u \in \mathcal{R}}^\top,$$
$$\hat{Z}_o = (\dot{o}\hat{\mathbf{p}}_u)_{u \in \hat{\mathcal{R}}}^\top.$$

In addition, we denote by $A$, the horizontal concatenation of the $A_o$ for all $o \in \Sigma$. Thus, we have

$$A = \begin{pmatrix} A_{o_1} \dots A_{o_{|\Sigma|}} \end{pmatrix}.$$

The proof of Theorem 1 is decomposed in four parts. First, we bound with high probability the maximum error for all $u \in \mathcal{P}$ between $\mathbf{d}_u$ and $\hat{\mathbf{d}}_u$ in norm $\ell_2$. Here, we use concentration inequalities like in (Hsu et al., 2012). Then, this error is propagated through the SPA using (Gillis & Vavasis, 2014). This allows bounding the perturbations in $\hat{Z}$ and $\hat{Z}_o$. Next, we analyze how solutions of the quadratic programming problems are perturbed using (Lőtstedt, 1983). Finally, perturbations in the estimated parameters of the PFA are multiplied together in a non-trivial way to finish the proof.

---
[1] now with Google DeepMind

## 1.2. Sampling errors

Let be $N$ independent sequences drawn from the target distribution $p$. These sequences are used to build empirical estimates $\hat{\mathbf{d}}_u$ of $\mathbf{d}_u$ for all $u$ in $\mathcal{P}$. The first step consists in bounding with high probability the error $\epsilon^{\text{est}} = \max_{u \in \mathcal{P}} \left\| \mathbf{d}_u - \hat{\mathbf{d}}_u \right\|_2$. We start by recalling a result in (**?**)Proposition 19]hsu2012spectral that uses the McDiarmid inequality (McDiarmid, 1989).

For all $u \in \mathcal{P}$, let $\mathbf{d}_u^\infty$ be an infinite vector such that $\forall v \in \Sigma^\star$, $\mathbf{d}_u^\infty[v] = \mathbb{P}\left(z = v | u\right) = \frac{p(uv)}{p(u\Sigma^\star)}$, where $z$ is a random variable with value in $\Sigma^\star$ drawn from $p_u$. Let $\hat{\mathbf{d}}_u^\infty$ be an estimator of $\mathbf{d}_u^\infty$ built from $n_u$ i.i.d. copies of $z$ denoted $z_i$.

**Lemma 1.** *We the previous notations, for all $\delta_u > 0$ we have*

$$\mathbb{P}\left( \left\| \mathbf{d}_u^\infty - \hat{\mathbf{d}}_u^\infty \right\|_2 \geq \frac{1}{\sqrt{n_u}} \left( 1 + \sqrt{\log\left(\frac{1}{\delta_u}\right)} \right) \right) \leq \delta_u.$$

In the sequel, we note $n = \min_{u \in \mathcal{P}} n_u$. In particular, $n_\varepsilon = N$.

**Proposition 1.** *With the previous notation, for all $\delta \in [0,1]$ we have*

$$\mathbb{P}\left( \epsilon^{est} \leq \frac{1}{\sqrt{n}} \left( 1 + \sqrt{\log\left(\frac{|\mathcal{P}|}{\delta}\right)} \right) \right) \geq 1 - \delta.$$

Before proving Proposition 1, we make few remarks. As $\mathbf{d}_u$ stands for a conditional distribution, the bound on $\epsilon^{\text{est}}$ depends necessarily on $n$. In addition, the bound depends on $\log(|\mathcal{P}|)$. We could obtain a bound independent of the dimension ($|\mathcal{P}|$) using (Denis et al., 2014) but the bound would be much more complicated. So, we kept a dimension dependent results.

*Proof.* By Lemma 1, we have

$$\mathbb{P}\left( \max_{u \in \mathcal{P}} \left\| \mathbf{d}_u^\infty - \hat{\mathbf{d}}_u^\infty \right\|_2 \leq \frac{1}{\sqrt{n}} \left( 1 + \sqrt{\log\left(\frac{|\mathcal{P}|}{\delta}\right)} \right) \right)$$

$$= 1 - \mathbb{P}\left( \exists u \in \mathcal{P}, \left\| \mathbf{d}_u^\infty - \hat{\mathbf{d}}_u^\infty \right\|_2 \geq \frac{1}{\sqrt{n}} \left( 1 + \sqrt{\log\left(\frac{|\mathcal{P}|}{\delta}\right)} \right) \right)$$

$$\geq 1 - \sum_{u \in \mathcal{P}} \mathbb{P}\left( \left\| \mathbf{d}_u^\infty - \hat{\mathbf{d}}_u^\infty \right\|_2 \geq \frac{1}{\sqrt{n}} \left( 1 + \sqrt{\log\left(\frac{|\mathcal{P}|}{\delta}\right)} \right) \right)$$

$$\geq 1 - \sum_{u \in \mathcal{P}} \mathbb{P}\left( \left\| \mathbf{d}_u^\infty - \hat{\mathbf{d}}_u^\infty \right\|_2 \geq \frac{1}{\sqrt{n_u}} \left( 1 + \sqrt{\log\left(\frac{|\mathcal{P}|}{\delta}\right)} \right) \right)$$

$$\geq 1 - \sum_{u \in \mathcal{P}} \frac{\delta}{|\mathcal{P}|} = 1 - \delta.$$

Next, by the definition of the norm $\ell_2$, we have $\left\| \mathbf{d}_u - \hat{\mathbf{d}}_u \right\|_2 \leq 2 \left\| \mathbf{d}_u^\infty - \hat{\mathbf{d}}_u^\infty \right\|_2$ because some coordinate of $\mathbf{d}_u^\infty - \hat{\mathbf{d}}_u^\infty$ can appear twice in $\mathbf{d}_u - \hat{\mathbf{d}}_u$. Finally, taking the maximum on $u \in \mathcal{P}$ leads to the result. $\qquad \square$

## 1.3. Perturbations in the convex hull

In this Section, we focus on how estimation errors perturb the identification of the convex hull. The analysis mainly use the results in (Gillis & Vavasis, 2014) on the robustness of SPA. First, we recall that we assumed $\text{rang}((\mathbf{d}_u)_{u \in \mathcal{R}}) = d$ as required by SPA. As the basis is complete, we also have that $\text{rang}(Z) = d$. In addition, $\sigma_d$ is the lowest positive singular value of $Z$ and so of $(\mathbf{d}_u)_{u \in \mathcal{R}}$ too. We denote $K = \max_{u \in \mathcal{R}} \|\mathbf{d}_u\|_2$.

As any reordering of $\mathcal{R}$ is inconsequential in CH-PRFA, we can assume without loss of generality that the $i^{\text{th}}$ element of $\mathcal{R}$, denoted $\mathcal{R}(i)$, correspond to $\hat{\mathcal{R}}(i)$. So, for clarity we slightly abuse the notation by denoting with the subscript $u$, both

$\mathcal{R}(i)$ and $\hat{\mathcal{R}}(i)$, even if they differ. For example, $\max_{i \in [1,d]} \left\| \hat{\mathbf{d}}_{\hat{\mathcal{R}}(i)} - \mathbf{d}_{\mathcal{R}(i)} \right\|_2$ becomes $\max_{u \in \hat{\mathcal{R}}} \left\| \hat{\mathbf{d}}_u - \mathbf{d}_u \right\|_2$. Similarly, $\max_{i \in [1,d]} \left| \frac{\hat{p}(\hat{\mathcal{R}}(i))}{\hat{p}(\hat{\mathcal{R}}(i)\Sigma^\star)} - \frac{p(\mathcal{R}(i))}{p(\mathcal{R}(i)\Sigma^\star)} \right|$ becomes $\max_{u \in \hat{\mathcal{R}}} \left| \frac{\hat{p}(u)}{\hat{p}(u\Sigma^\star)} - \frac{p(u)}{p(u\Sigma^\star)} \right|$. Now, we define the perturbations in the convex hull to be $\epsilon^{\text{conv}} = \max_{u \in \hat{\mathcal{R}}} \left\| \hat{\mathbf{d}}_u - \mathbf{d}_u \right\|_2$.

**Proposition 2.** *With the previous notations, if $\epsilon^{est} \leq \frac{\sigma_d^3}{648\sqrt{d}}$ then,*

$$\epsilon^{conv} \leq 162 \frac{\epsilon^{est}}{\sigma_d^2}.$$

*Proof.* The proof is a direct application of the Theorem 3 in (Gillis & Vavasis, 2014) followed by some simplifications. This Theorem shows that if

$$\max_{u \in \mathcal{P}} \left\| \hat{\mathbf{d}}_u - \mathbf{d}_u \right\|_2 < \sigma_d \min \left( \frac{1}{2\sqrt{d-1}}, \frac{1}{4} \right) \left( 1 + 80 \frac{K^2}{\sigma_d^2} \right)^{-1},$$

then

$$\max_{u \in \hat{\mathcal{R}}} \left\| \hat{\mathbf{d}}_u - \mathbf{d}_u \right\|_2 < \max_{u \in \mathcal{P}} \left\| \hat{\mathbf{d}}_u - \mathbf{d}_u \right\|_2 \left( 1 + 80 \frac{K^2}{\sigma_d^2} \right).$$

First, we have

$$\frac{1}{4\sqrt{d}} \leq \min \left( \frac{1}{2\sqrt{d-1}}, \frac{1}{4} \right).$$

Next, as $\mathbf{d}_u$ contains probabilities, $\|\mathbf{d}_u\|_2 \leq \|\mathbf{d}_u\|_1 \leq 1$ and so, we have $K \leq 1$. Using Lemme 4 in (Gillis & Vavasis, 2014), we also have that $\sigma_d \leq K$. This gives the following inequality

$$\left( 1 + 80 \frac{K^2}{\sigma_d^2} \right) \leq \frac{1}{\sigma_d^2} \left( \sigma_d^2 + 80K^2 \right) \leq \frac{81}{\sigma_d^2}.$$

Thus, we can simplify the two bounds as follows,

$$\sigma_d \min \left( \frac{1}{2\sqrt{d-1}}, \frac{1}{4} \right) \left( 1 + 80 \frac{K^2}{\sigma_d^2} \right)^{-1} \geq \sigma_d \frac{1}{4\sqrt{d}} \left( 1 + 80 \frac{K^2}{\sigma_d^2} \right)^{-1} \geq \frac{\sigma_d^3}{324\sqrt{d}},$$

and,

$$\max_{u \in \mathcal{P}} \left\| \hat{\mathbf{d}}_u - \mathbf{d}_u \right\|_2 \left( 1 + 80 \frac{K^2}{\sigma_d^2} \right) \leq 81 \frac{\max_{u \in \mathcal{P}} \left\| \hat{\mathbf{d}}_u - \mathbf{d}_u \right\|_2}{\sigma_d^2}.$$

Finally, we conclude using

$$\max_{u \in \mathcal{P}} \left\| \hat{\mathbf{d}}_u - \mathbf{d}_u \right\|_2 \leq 2\epsilon^{\text{est}}.$$

$\square$

## 1.4. Perturbations in solutions of the quadratic programming problems

In this Section, we first introduce a general form of the quadratic programming problems involved in CH-PRFA. This form allows us treating simultaneously all the minimization problems together up to one point. Then, in the next Section we will focus on each particular problems. The general form we consider is the following

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} \|Q\mathbf{x} + \mathbf{q}\|_2 \tag{1}$$

$$\text{s.t.} \begin{cases} B\mathbf{x} + \mathbf{b} \geq 0, \\ C\mathbf{x} + \mathbf{c} = 0 \end{cases}. \tag{2}$$

Here, we are interested in bounding $\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2$ by a function of $\left\| \hat{Q} - Q \right\|_2, \|\hat{\mathbf{q}} - \mathbf{q}\|_2, \left\| \hat{B} - B \right\|_2, \left\| \hat{\mathbf{b}} - \mathbf{b} \right\|_2, \left\| \hat{C} - C \right\|_2, \|\hat{\mathbf{c}} - \mathbf{c}\|_2$, where $\hat{\mathbf{x}}^*$ is the solution to the perturbed problem.

Next, we show that the minimization problems in CH-PRFA can be written in the general form of Equation (1). Let $\hat{\mathbf{x}} = \hat{\boldsymbol{\alpha}}_0 = (\hat{a}_\varepsilon^u)_{u \in \hat{\mathcal{R}}}$, then we have

$$\hat{Q} = \hat{Z}^\top, \qquad\qquad \hat{B} = I, \qquad\qquad \hat{C} = \mathbf{1}^\top,$$
$$\hat{\mathbf{q}} = -\hat{H}_{\mathcal{B}}^\top \mathbf{1}_\varepsilon, \qquad\qquad \hat{\mathbf{b}} = 0, \qquad\qquad \hat{\mathbf{c}} = -1.$$

Similarly, for all $u \in \hat{\mathcal{R}}$, let $\hat{\mathbf{x}}_u$ be a vector of size $d\,|\Sigma|$ such that

$$\hat{\mathbf{x}}_u^\top = \underbrace{\left( \begin{array}{ccc} \hat{A}_{o_1}[u,:] & \cdots & \hat{A}_{o_{|\Sigma|}}[u,:] \end{array} \right)}_{d|\Sigma| \text{ columns}} = \hat{A}[u,:].$$

then we have,

$$\hat{Q} = \underbrace{\begin{pmatrix} \hat{Z}^\top & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \hat{Z}^\top \end{pmatrix}}_{d|\Sigma| \text{ columns}},$$

$$\hat{\mathbf{q}}_u^\top = - \underbrace{\left( \begin{array}{ccc} \hat{Z}_1[u,:] & \cdots & \hat{Z}_{|\Sigma|}[u,:] \end{array} \right)}_{d|\Sigma| \text{ columns}},$$

$$\hat{B} = I,$$
$$\hat{\mathbf{b}} = 0,$$
$$\hat{C} = \mathbf{1}^\top,$$
$$\hat{\mathbf{c}}_u = \frac{\hat{p}(u)}{\hat{p}(u\Sigma^\star)} - 1.$$

### 1.4.1. EXISTENCE AND UNIQUENESS OF THE SOLUTION

In this Section, we check that the solutions $\mathbf{x}^*$ and $\hat{\mathbf{x}}^*$ to the previous unperturbed and perturbed quadratic programming problems exists. We denote the kernel of a matrix $M$ by $\mathrm{N}(M)$ and the range by $\mathrm{R}(M)$. We denote by $E$ the matrix and $\mathbf{e}$ the vector such that

$$E = \left( \begin{array}{c} B \\ C \end{array} \right), \qquad\qquad \mathbf{e} = \left( \begin{array}{c} \mathbf{b} \\ \mathbf{c} \end{array} \right).$$

First, its straightforward to verify that the sets delimited by the linear constraints are not empty for both the unperturbed and perturbed problems. Thus, using Theorem 1 in (Lőtstedt, 1983), we know that $\mathbf{x}^*$ et $\hat{\mathbf{x}}^*$ exists.

Next, we show that

$$\mathrm{N}(Q) = \mathrm{N}(\hat{Q}) = \{0\}, \tag{3}$$
$$\mathrm{N}(E) = \mathrm{N}(\hat{E}) = \{0\}. \tag{4}$$

For all the quadratic programming problems involved in CH-PRFA, we have that $B = \hat{B} = I$. So $\mathrm{N}(E) = \mathrm{N}(\hat{E}) = \{0\}$. Moreover, as $\mathrm{rang}(Z) = d$, we also have that $\mathrm{N}(Q) = \{0\}$. Lastly, we have that $\mathrm{N}(\hat{Q}) = \{0\}$, with probability 1, by density of invertible matrices. Again using Theorem 1 in (Lőtstedt, 1983), we know that $\hat{\mathbf{x}}^*$ and $\mathbf{x}^*$ are unique.

### 1.4.2. BOUND ON THE PERTURBATIONS

First, we verify that for all quadratic programming problems in CH-PRFA, we have

$$B = \hat{B}, \qquad\qquad \mathbf{b} = \hat{\mathbf{b}}, \qquad\qquad C = \hat{C}.$$

This implies $E = \hat{E}$. Next, we adopt the same notations than in (Lőtstedt, 1983) by denoting $G = EQ^\dagger$. So, we have that $\hat{G} = E\hat{Q}^\dagger$. In addition, we denote $\mathbf{v}^* = 2(G^\dagger)^\top(Q\mathbf{x}^* + P_{\mathrm{R}(Q)}\mathbf{q})$ where $P_{\mathrm{R}(Q)}$ is a projector on $\mathrm{R}(Q)$. As $(\boldsymbol{\alpha}_0, A, \boldsymbol{\alpha}_\infty)$ defines a PRFA by hypothesis, we have exactly that $Q\mathbf{x}^* + P_{\mathrm{R}(Q)}\mathbf{q} = 0$. Thus, we have $\mathbf{v}^* = 0$. As the basis is complete, we have $P_{\mathrm{R}(Q)}\mathbf{q} = \mathbf{q}$. These particular properties will allow us simplifying the results of Theorem 3 in (Lőtstedt, 1983).

**Theorem 2.** *If Equation* (3) *holds, let* $\mathbf{x}^*$ *be the solution of* (1) *and* $\hat{\mathbf{x}}^*$ *be a perturbed solution. If* $Q\mathbf{x}^* + P_{\mathrm{R}(Q)}\mathbf{q} = 0$, $P_{\mathrm{R}(Q)}\mathbf{q} = \mathbf{q}$, $B = \hat{B}$, $\mathbf{b} = \hat{\mathbf{b}}$ *and* $C = \hat{C}$, *then*

$$\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2 \le 2\|\mathbf{q}\|_2 \left(1 + \left\|\hat{Q}\right\|_2\left\|\hat{Q}^\dagger\right\|_2\|E\|_2\|E^\dagger\|_2\right)\left\|\hat{Q}^\dagger - Q^\dagger\right\|_2 + \left\|\hat{Q}\right\|_2\left\|\hat{Q}^\dagger\right\|_2\|E^\dagger\|_2\|\hat{\mathbf{e}} - \mathbf{e}\|_2 + 2\left\|\hat{Q}^\dagger\right\|_2\|\hat{\mathbf{q}} - \mathbf{q}\|_2. \tag{5}$$

*Proof.* We start by the unsimplified bound of the Theorem 3 in (Lőtstedt, 1983),

$$\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2 \le \left\|Q\mathbf{x}^* - P_{\mathrm{N}(Q^\top)}\mathbf{q}\right\|_2 \left(\left\|\hat{Q}^\dagger - Q^\dagger\right\|_2 + \left\|\hat{Q}^\dagger\right\|_2\left\|\hat{G}^\dagger\right\|_2\left\|\hat{G} - G\right\|_2\right)$$
$$+ \left\|\hat{Q}^\dagger\right\|_2\left(2\left\|\frac{1}{2}\left(\hat{G} - G\right)^\top\mathbf{v}^* - (\hat{\mathbf{q}} - \mathbf{q})\right\|_2 + \left\|\hat{G}^\dagger\right\|_2\|\hat{\mathbf{e}} - \mathbf{e}\|_2\right).$$

As $P_{\mathrm{N}(Q^\top)} = P_{\mathrm{R}(Q)}$ and $Q\mathbf{x}^* + P_{\mathrm{R}(Q)}\mathbf{q} = 0$ then

$$\left\|Q\mathbf{x}^* - P_{\mathrm{N}(Q^\top)}\mathbf{q}\right\|_2 = 2\|\mathbf{q}\|_2.$$

In addition, as $\mathbf{v}^* = 0$ we have

$$\left\|\frac{1}{2}\left(\hat{G} - G\right)^\top\mathbf{v}^* - (\hat{\mathbf{q}} - \mathbf{q})\right\|_2 = \|\hat{\mathbf{q}} - \mathbf{q}\|_2.$$

Moreover, as $G = EQ^\dagger$ et $E = \hat{E}$, we obtain

$$\left\|\hat{G} - G\right\|_2 = \left\|E\left(\hat{Q}^\dagger - Q^\dagger\right)\right\|_2 \le \|E\|_2\left\|\hat{Q}^\dagger - Q^\dagger\right\|_2.$$

In addition, as $\mathrm{N}(E) = \{0\}$ and $\mathrm{N}(Q) = \{0\}$, we have $\hat{G}^\dagger = \left(E\hat{Q}^\dagger\right)^\dagger = \hat{Q}E^\dagger$ and using the triangle inequality

$$\left\|\hat{G}^\dagger\right\|_2 \le \|E^\dagger\|_2\left\|\hat{Q}\right\|_2.$$

Finally, we obtain

$$\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2 \le 2\|\mathbf{q}\|_2 \left(1 + \left\|\hat{Q}\right\|_2\left\|\hat{Q}^\dagger\right\|_2\|E\|_2\|E^\dagger\|_2\right)\left\|\hat{Q}^\dagger - Q^\dagger\right\|_2 + \left\|\hat{Q}\right\|_2\left\|\hat{Q}^\dagger\right\|_2\|E^\dagger\|_2\|\hat{\mathbf{e}} - \mathbf{e}\|_2 + 2\left\|\hat{Q}^\dagger\right\|_2\|\hat{\mathbf{q}} - \mathbf{q}\|_2.$$
$\square$

To further simplify, we need a Lemma in (Wedin, 1972) that bounds the perturbation in the Moore-Pseudo inverse.

**Lemma 2.** *If* $\left\|\hat{Q} - Q\right\|_2\|Q^\dagger\|_2 \le \kappa < 1$ *then*

$$\left\|\hat{Q}^\dagger\right\|_2 \le \frac{1}{1-\kappa}\|Q^\dagger\|_2,$$

*and*

$$\left\|\hat{Q}^\dagger - Q^\dagger\right\|_2 \le \sqrt{2}\left\|\hat{Q}^\dagger\right\|_2\|Q^\dagger\|_2\left\|\hat{Q} - Q\right\|_2 \le \frac{\sqrt{2}}{1-\kappa}\|Q^\dagger\|_2^2\left\|\hat{Q} - Q\right\|_2.$$

Thus, Lemma 2 implies that

$$\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2 \le \frac{2\sqrt{2}}{1-\kappa}\|\mathbf{q}\|_2 \left(1 + \frac{1}{1-\kappa}\left\|\hat{Q}\right\|_2\|Q^\dagger\|_2\|E\|_2\|E^\dagger\|_2\right)\|Q^\dagger\|_2^2\left\|\hat{Q} - Q\right\|_2$$
$$+ \frac{1}{1-\kappa}\left\|\hat{Q}\right\|_2\|Q^\dagger\|_2\|E^\dagger\|_2\|\hat{\mathbf{e}} - \mathbf{e}\|_2 \tag{6}$$
$$+ \frac{2}{1-\kappa}\|Q^\dagger\|_2\|\hat{\mathbf{q}} - \mathbf{q}\|_2.$$

In the sequel, we analyze each quadratic programming problems separately.

## 1.5. Perturbations in the automata parameters

### 1.5.1. PERTURBATIONS IN THE INITIAL PROBABILITIES

In this section, to propose a bound on the perturbation in $\hat{\alpha}_0$, we focus on the following problem,

$$\{\hat{a}_\varepsilon^u\} = \underset{\{a_\varepsilon^u\}}{\operatorname{argmin}} \left\| \hat{\mathbf{p}} - \sum_{u \in \hat{\mathcal{R}}} a_\varepsilon^u \hat{\mathbf{p}}_u \right\|_2 \quad \text{s.t.} \sum_{u \in \hat{\mathcal{R}}} a_\varepsilon^u = 1 \text{ and } a_\varepsilon^u \geq 0.$$

**Proposition 3.** *Let $\kappa$ be a real in $[0, 1[$, if $\epsilon^{conv} \leq \frac{1}{\sqrt{d}}\kappa\sigma_d$, then*

$$\|\hat{\alpha}_0 - \alpha_0\|_1 \leq 2(\sqrt{2} + 2)\frac{d^2}{\sigma_d^3 (1 - \kappa)^2}\epsilon^{conv} + 2\frac{\sqrt{d}}{\sigma_d (1 - \kappa)}\epsilon^{est}.$$

*Proof.* First, we recall that $\hat{\mathbf{e}} = \mathbf{e} = \begin{pmatrix} \mathbf{b}^\top & \mathbf{c}^\top \end{pmatrix}^\top = \begin{pmatrix} 0 & \ldots & 0 & -1 \end{pmatrix}^\top$. By replacing in Equation (6), we obtain that

$$\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2 \leq \frac{2\sqrt{2}}{1 - \kappa}\|\mathbf{q}\|_2 \left(1 + \frac{1}{1 - \kappa}\|\hat{Q}\|_2\|Q^\dagger\|_2\|E\|_2\|E^\dagger\|_2\right)\|Q^\dagger\|_2^2\|\hat{Q} - Q\|_2 + \frac{2}{1 - \kappa}\|Q^\dagger\|_2\|\hat{\mathbf{q}} - \mathbf{q}\|_2.$$

Next, we have

$$E^\top E = B^\top B + C^\top C = I_d + U_d,$$

where $I_d$ is the identity matrix of dimension $d$ and $U_d$ is the unity matrix (all the coefficients equal 1) of dimension $d$. We denote by $\lambda_{\max}(M)$ (resp. $\lambda_{\min}(M)$) the largest (resp. smallest) eigenvalue of $M$. Note that, the eigenvalues of $U_d$ are $d$ with multiplicity one and 0 with multiplicity $d - 1$. This implies that

$$\|E\|_2^2 = \lambda_{\max}(E^\top E) = \lambda_{\max}(I_d + U_d) = 1 + d,$$

and

$$\|E^\dagger\|_2^2 = \lambda_{\min}(E^\top E) = \lambda_{\min}(I_d + U_d) = 1.$$

These properties allow us to simplify again Equation (6),

$$\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2 \leq \frac{2\sqrt{2}}{1 - \kappa}\|\mathbf{q}\|_2 \left(1 + \frac{\sqrt{1 + d}}{1 - \kappa}\|\hat{Q}\|_2\|Q^\dagger\|_2\right)\|Q^\dagger\|_2^2\|\hat{Q} - Q\|_2 + \frac{2}{1 - \kappa}\|Q^\dagger\|_2\|\hat{\mathbf{q}} - \mathbf{q}\|_2.$$

In addition, as $Q = Z^\top$, we have that $\|Q^\dagger\|_2 = \frac{1}{\sigma_d}$, where $\sigma_d$ is the smallest eigenvalue of $Z$.

Now, we focus on $\|\mathbf{q}\|_2$ and $\|\hat{Q}\|_2$, where $\mathbf{q} = -H_{\mathcal{B}}^\top \mathbf{1}_\varepsilon$ and $\hat{Q} = \hat{Z}^\top$. First, as $\mathbf{q}$ contains probabilities, we have that

$$\|\mathbf{q}\|_2 \leq \|\mathbf{q}\|_1 \leq \|H_{\mathcal{B}}^\top \mathbf{1}_\varepsilon\|_1 = \sum_{u \in \mathcal{S}} p(u) \leq 1.$$

Secondly, the Hölder inequality implies

$$\|\hat{Q}\|_2 \leq \sqrt{\|\hat{Q}\|_1 \|\hat{Q}\|_\infty}.$$

On one hand, we have $\|\hat{Q}\|_1 = \|\hat{Z}^\top\|_1 = \max_{u \in \hat{\mathcal{R}}} \sum_{v \in \mathcal{S}} \frac{\hat{p}(uv)}{\hat{p}(u\Sigma)} \leq 1$. On the other hand, we show that $\|\hat{Q}\|_\infty = \|\hat{Z}^\top\|_\infty = \max_{v \in \mathcal{S}} \sum_{u \in \hat{\mathcal{R}}} \frac{\hat{p}(uv)}{\hat{p}(u\Sigma)} \leq |\hat{\mathcal{R}}| = d$. Thus, we have

$$\|\hat{Q}\|_2 \leq \sqrt{d}.$$

Finally, we obtain that

$$\|\hat{\alpha}_0 - \alpha_0\|_2 \leq \frac{2\sqrt{2}}{\sigma_d^2 (1 - \kappa)} \left(1 + \frac{\sqrt{d(1 + d)}}{\sigma_d (1 - \kappa)}\right)\|\hat{Z} - Z\|_2 + \frac{2}{\sigma_d (1 - \kappa)}\|\hat{\mathbf{q}} - \mathbf{q}\|_2.$$

When we considered the perturbation in the convex hull, we showed that $\sigma_d \leq 1$. As $\kappa < 1$, we have

$$1 + \frac{\sqrt{d(1+d)}}{\sigma_d(1-\kappa)} \leq \frac{1 + \sqrt{d(1+d)}}{\sigma_d(1-\kappa)} \leq \frac{(1+\sqrt{2})d}{\sigma_d(1-\kappa)},$$

because for $d \geq 1$, we have $1 + \sqrt{d(1+d)} \leq (1+\sqrt{2})d$. Combining the previous inequalities leads to

$$\|\hat{\boldsymbol{\alpha}}_0 - \boldsymbol{\alpha}_0\|_2 \leq \frac{2(\sqrt{2}+2)d}{\sigma_d^3(1-\kappa)^2}\left\|\hat{Z} - Z\right\|_2 + \frac{2}{\sigma_d(1-\kappa)}\|\hat{\mathbf{q}} - \mathbf{q}\|_2.$$

The next step is to insert the bounds on $\left\|\hat{Z} - Z\right\|_2$ and $\|\hat{\mathbf{q}} - \mathbf{q}\|_2$. First, we have

$$\|\hat{\mathbf{q}} - \mathbf{q}\|_2 = \left\|\left(\hat{H}_\mathcal{B} - H_\mathcal{B}\right)^\top \mathbf{1}_\varepsilon\right\|_2 = \sqrt{\sum_{v \in \mathcal{S}}(\hat{p}(v) - p(v))^2} \leq \sqrt{\sum_{v \in \Sigma^\star}(\hat{p}(v) - p(v))^2} = \left\|\hat{\mathbf{d}}_\varepsilon^\infty - \mathbf{d}_\varepsilon^\infty\right\|_2.$$

Next, using properties of the norm $\ell_2$ and of the Frobenius norm, we show that

$$\left\|\hat{Z} - Z\right\|_2 \leq \left\|\hat{Z} - Z\right\|_F = \sqrt{\sum_{u \in \hat{\mathcal{R}}}\left(\left\|\hat{Z}[u,:] - Z[u,:]\right\|_2\right)^2} \leq \sqrt{\sum_{u \in \hat{\mathcal{R}}}\left(\left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2\right)^2} \leq \sqrt{d}\max_{u \in \hat{\mathcal{R}}}\left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2.$$

Thus, as

$$\left\|\hat{Q} - Q\right\|_2 = \left\|\hat{Z} - Z\right\|_2 \leq \sqrt{d}\max_{u \in \hat{\mathcal{R}}}\left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2 \kappa\sigma_d,$$

we have that $\left\|\hat{Q} - Q\right\|_2 \leq \kappa\sigma_d$ holds if $\max_{u \in \hat{\mathcal{R}}}\left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2 \leq \frac{1}{\sqrt{d}}\kappa\sigma_d$. Finally, using norms properties, we obtain $\|\hat{\boldsymbol{\alpha}}_0 - \boldsymbol{\alpha}_0\|_1 \leq \sqrt{d}\|\hat{\boldsymbol{\alpha}}_0 - \boldsymbol{\alpha}_0\|_2$ and

$$\|\hat{\boldsymbol{\alpha}}_0 - \boldsymbol{\alpha}_0\|_1 \leq 2(\sqrt{2}+2)\frac{d^2}{\sigma_d^3(1-\kappa)^2}\max_{u \in \hat{\mathcal{R}}}\left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2 + 2\frac{\sqrt{d}}{\sigma_d(1-\kappa)}\left\|\hat{\mathbf{d}}_\varepsilon^\infty - \mathbf{d}_\varepsilon^\infty\right\|_2.$$

To conclude, just substitute in the final inequality

$$\max_{u \in \hat{\mathcal{R}}}\left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2 = \epsilon^{\text{conv}}$$

and

$$\left\|\hat{\mathbf{d}}_\varepsilon^\infty - \mathbf{d}_\varepsilon^\infty\right\|_2 \leq \max_{u \in \mathcal{P}}\left\|\hat{\mathbf{d}}_u^\infty - \mathbf{d}_u^\infty\right\|_2 = \epsilon^{\text{est}}.$$

$\square$

### 1.5.2. Perturbations on transition probabilities

In this Section, we focus on the following problems, for all $u \in \hat{\mathcal{R}}$

$$\{\hat{a}_{u,o}^v\} = \underset{\{a_{u,o}^v\}}{\text{argmin}}\sum_{o \in \Sigma}\left\|\dot{o}\hat{\mathbf{p}}_u - \sum_{v \in \hat{\mathcal{R}}}a_{u,o}^v\hat{\mathbf{p}}_v\right\|_2 \quad \text{s.t.} \sum_{v \in \hat{\mathcal{R}}, o \in \Sigma}a_{u,o}^v\dot{o} = 1 - \hat{\mathbf{p}}_u\mathbf{1}_\epsilon \text{ and } a_{u,o}^v \geq 0.$$

**Proposition 4.** *Let $\kappa$ be a real in $[0,1[$, if $\epsilon^{\text{conv}} \leq \frac{1}{\sqrt{d}}\kappa\sigma_d$, then*

$$\left\|\hat{A} - A\right\|_\infty \leq 2(\sqrt{2}+2)\frac{d^2\sqrt{|\Sigma|}}{\sigma_d^3(1-\kappa)^2}\epsilon^{\text{conv}} + \frac{d}{\sigma_d(1-\kappa)}\|\hat{\boldsymbol{\alpha}}_\infty - \boldsymbol{\alpha}_\infty\|_\infty + 2\frac{\sqrt{d}}{\sigma_d(1-\kappa)}\epsilon^{\text{conv}}.$$

The proof of Proposition 4 follows the same steps than Proposition 3 for each individual problems and then combine the results for all $u \in \hat{\mathcal{R}}$.

*Proof.* First, we have

$$E^\top E = B^\top B + C^\top C$$
$$= I_{d|\Sigma|} + U_{d|\Sigma|},$$

This implies that

$$\|E\|_2^2 = \lambda_{\max}\left(E^\top E\right) = \lambda_{\max}\left(I_{d|\Sigma|} + U_{d|\Sigma|}\right) = 1 + d\,|\Sigma|\,.$$

In addition, we have

$$\left\|E^\dagger\right\|_2^2 = \lambda_{\min}\left(E^\top E\right) = \lambda_{\min}\left(I_{d|\Sigma|} + U_{d|\Sigma|}\right) = 1.$$

By replacing in Equation (6), for the problem associated with $u \in \hat{\mathcal{R}}$, and assuming $\left\|\hat{Q} - Q\right\|_2 \leq \sigma_d\kappa$ we obtain,

$$\|\hat{\mathbf{x}}_u^* - \mathbf{x}_u^*\|_2 \leq \frac{2\sqrt{2}}{1-\kappa}\|\mathbf{q}_u\|_2 \left(1 + \frac{\sqrt{1 + d\,|\Sigma|}}{1-\kappa}\left\|\hat{Q}\right\|_2\|Q^\dagger\|_2\right)\|Q^\dagger\|_2^2\left\|\hat{Q} - Q\right\|_2$$
$$+ \frac{1}{1-\kappa}\left\|\hat{Q}\right\|_2\|Q^\dagger\|_2\|\hat{\mathbf{e}}_u - \mathbf{e}_u\|_2$$
$$+ \frac{2}{1-\kappa}\|Q^\dagger\|_2\|\hat{\mathbf{q}}_u - \mathbf{q}_u\|_2.$$

Moreover, as $Q$ is block diagonal, $Q$ has the same eigenvalues than $Z^\top$ and $\|Q^\dagger\|_2 = \|Z^\top\|_2 = \frac{1}{\sigma_d}$.

Now, we analyze $\|\mathbf{q}_u\|_2$ and $\left\|\hat{Q}\right\|_2$. First, we have that

$$\|\mathbf{q}_u\|_2 \leq \|\mathbf{q}_u\|_1 \leq \sum_{o\in\Sigma}\sum_{v\in\mathcal{S}}\frac{p(uov)}{p(u\Sigma\Sigma^\star)} \leq \frac{p(u\Sigma^\star)}{p(u\Sigma^\star)} = 1.$$

Next, the Hlder inequality implies that

$$\left\|\hat{Q}\right\|_2 = \left\|\hat{Z}^\top\right\|_2 \leq \sqrt{\left\|\hat{Z}^\top\right\|_1\left\|\hat{Z}^\top\right\|_\infty}.$$

On one hand, we have that $\left\|\hat{Z}^\top\right\|_1 = \max_{u\in\hat{\mathcal{R}}}\sum_{v\in\mathcal{S}}\frac{\hat{p}(uv)}{\hat{p}(u\Sigma)} \leq 1$. On the other hand, we show that $\left\|\hat{Z}^\top\right\|_\infty = \max_{v\in\mathcal{S}}\sum_{u\in\hat{\mathcal{R}}}\frac{\hat{p}(uv)}{\hat{p}(u\Sigma)} \leq \left|\hat{\mathcal{R}}\right| = d$. Thus, we have

$$\left\|\hat{Q}\right\|_2 \leq \sqrt{d}.$$

The two previous inequalities implies that,

$$\|\hat{\mathbf{x}}_u^* - \mathbf{x}_u^*\|_2 \leq \frac{2\sqrt{2}}{\sigma_d^2(1-\kappa)}\left(1 + \frac{\sqrt{d(1 + d\,|\Sigma|)}}{\sigma_d(1-\kappa)}\right)\left\|\hat{Z} - Z\right\|_2 + \frac{\sqrt{d}}{\sigma_d(1-\kappa)}\|\hat{\mathbf{e}}_u - \mathbf{e}_u\|_2 + \frac{2}{\sigma_d(1-\kappa)}\|\hat{\mathbf{q}}_u - \mathbf{q}_u\|_2.$$

When we considered the perturbation in the convex hull, we showed that $\sigma_d \leq 1$. As $\kappa < 1$, we have

$$1 + \frac{\sqrt{d\,(1 + d\,|\Sigma|)}}{\sigma_d\,(1-\kappa)} \leq \frac{1 + \sqrt{d\,(1 + d\,|\Sigma|)}}{\sigma_d\,(1-\kappa)} \leq \frac{(1+\sqrt{2})d\sqrt{|\Sigma|}}{\sigma_d\,(1-\kappa)},$$

because for $d \geq 1$, we have $1 + \sqrt{d\,(1+d)} \leq (1+\sqrt{2})d$. Replacing in the main inequality leads to

$$\|\hat{\mathbf{x}}_u^* - \mathbf{x}_u^*\|_2 \leq \frac{2(\sqrt{2}+2)d\sqrt{|\Sigma|}}{\sigma_d^3(1-\kappa)^2}\left\|\hat{Z} - Z\right\|_2 + \frac{\sqrt{d}}{\sigma_d(1-\kappa)}\|\hat{\mathbf{e}}_u - \mathbf{e}_u\|_2 + \frac{2}{\sigma_d(1-\kappa)}\|\hat{\mathbf{q}}_u - \mathbf{q}_u\|_2.$$

Because the norm $\ell_\infty$ is the maximum absolute row sum, we have

$$\left\|\hat{A} - A\right\|_\infty = \max_{u\in\hat{\mathcal{R}}}\left(\left\|\hat{A}[u,:] - A[u,:]\right)^\top\right\|_1 = \max_{u\in\hat{\mathcal{R}}}\|\hat{\mathbf{x}}_u - \mathbf{x}_u\|_1 \leq \sqrt{d}\max_{u\in\hat{\mathcal{R}}}\|\hat{\mathbf{x}}_u - \mathbf{x}_u\|_2.$$

The next step is to insert the bounds on $\left\|\hat{Z} - Z\right\|_2$ and $\left\|\hat{Z} - Z\right\|_2$, $\max_{u \in \hat{\mathcal{R}}} \|\hat{\mathbf{q}}_u - \mathbf{q}_u\|_2$ and $\max_{u \in \hat{\mathcal{R}}} \|\hat{\mathbf{e}}_u - \mathbf{e}_u\|_2$ First, using properties of the norm $\ell_2$ and of the Frobenius norm, we show that

$$\left\|\hat{Z} - Z\right\|_2 \le \left\|\hat{Z} - Z\right\|_F = \sqrt{\sum_{u \in \hat{\mathcal{R}}} \left(\left\|\hat{Z}[u,:] - Z[u,:]\right\|_2\right)^2} \le \sqrt{\sum_{u \in \hat{\mathcal{R}}} \left(\left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2\right)^2} \le \sqrt{d} \max_{u \in \hat{\mathcal{R}}} \left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2.$$

Thus, as

$$\left\|\hat{Q} - Q\right\|_2 = \left\|\hat{Z} - Z\right\|_2 \le \sqrt{d} \max_{u \in \hat{\mathcal{R}}} \left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2 \kappa\sigma_d,$$

the condition $\left\|\hat{Q} - Q\right\|_2 \le \kappa\sigma_d$ holds if, $\max_{u \in \hat{\mathcal{R}}} \left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2 \le \frac{1}{\sqrt{d}}\kappa\sigma_d$. Secondly, we have

$$\max_{u \in \hat{\mathcal{R}}} \|\hat{\mathbf{q}}_u - \mathbf{q}_u\|_2 = \max_{u \in \hat{\mathcal{R}}} \left\|\left( \hat{Z}_1[u,:] - Z_1[u,:] \quad \ldots \quad \hat{Z}_{|\Sigma|}[u,:] - Z_{|\Sigma|}[u,:] \right)\right\|_2$$
$$\le \max_{u \in \hat{\mathcal{R}}} \left\|\left( \hat{Z}[u,:] - Z[u,:] \quad \hat{Z}_1[u,:] - Z_1[u,:] \quad \ldots \quad \hat{Z}_{|\Sigma|}[u,:] - Z_{|\Sigma|}[u,:] \right)\right\|_2$$
$$\le \max_{u \in \hat{\mathcal{R}}} \left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2.$$

Finally, as $\hat{\mathbf{b}} = \mathbf{b}$, we have

$$\max_{u \in \hat{\mathcal{R}}} \|\hat{\mathbf{e}}_u - \mathbf{e}_u\|_2 = \max_{u \in \hat{\mathcal{R}}} \|\hat{\mathbf{c}}_u - \mathbf{c}_u\|_2 = \max_{u \in \hat{\mathcal{R}}} \left|\frac{\hat{p}(u)}{\hat{p}(u\Sigma^\star)} - \frac{p(u)}{p(u\Sigma^\star)}\right| = \|\hat{\boldsymbol{\alpha}}_\infty - \boldsymbol{\alpha}_\infty\|_\infty.$$

To conclude, we use the three last equalities and inequalities on $\max_{u \in \hat{\mathcal{R}}} \|\hat{\mathbf{q}}_u - \mathbf{q}_u\|_2$, $\left\|\hat{Q} - Q\right\|_2$ and $\max_{u \in \hat{\mathcal{R}}} \|\hat{\mathbf{e}}_u - \mathbf{e}_u\|_2$, in addition to $\max_{u \in \hat{\mathcal{R}}} \left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2 = \epsilon^{\mathrm{conv}}$ to replace in the main inequality. $\qquad\square$

### 1.5.3. PERTURBATIONS IN THE FINAL PROBABILITIES

The proof of Proposition 5 is much simpler than the previous as it does not involve quadratic programming problems.

**Proposition 5.**
$$\|\hat{\boldsymbol{\alpha}}_\infty - \boldsymbol{\alpha}_\infty\|_\infty \le \epsilon^{conv}.$$

*Proof.* By definition $\|\hat{\boldsymbol{\alpha}}_\infty - \boldsymbol{\alpha}_\infty\|_\infty = \max_{u \in \hat{\mathcal{R}}} \left|\frac{\hat{p}(\hat{\mathcal{R}}(i))}{\hat{p}(\hat{\mathcal{R}}(i)\Sigma^\star)} - \frac{p(u)}{p(u\Sigma^\star)}\right|$. Yet, we have

$$\max_{u \in \hat{\mathcal{R}}} \left|\frac{\hat{p}(u)}{\hat{p}(u\Sigma^\star)} - \frac{p(u)}{p(u\Sigma^\star)}\right| = \max_{u \in \hat{\mathcal{R}}} \sqrt{\left(\frac{\hat{p}(u)}{\hat{p}(u\Sigma^\star)} - \frac{p(u)}{p(u\Sigma^\star)}\right)^2} \le \max_{u \in \hat{\mathcal{R}}} \sqrt{\sum_{v \in \mathcal{S}} \left(\frac{\hat{p}(uv)}{\hat{p}(u\Sigma^\star)} - \frac{p(uv)}{p(u\Sigma^\star)}\right)^2} \le \max_{u \in \hat{\mathcal{R}}} \left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2.$$

To conclude, just substitute in the final inequality,

$$\left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2 = \epsilon^{\mathrm{conv}}.$$

$\qquad\square$

## 1.6. Perturbations in the distribution

Before proving Theorem 1, we prove the following Proposition from which Theorem 1 can be easily deduced.

**Proposition 6.** *Under the same hypothesis than Theorem 1, for all $0 < \delta < 1$, there exists a constant $K$ such that, for all $t > 0$, $\epsilon > 0$, with probability $1 - \delta$, if*

$$n \ge K \frac{t^2 d^4 |\Sigma|}{\epsilon^2 \sigma_d^{10}} \log\left(\frac{|\mathcal{P}|}{\delta}\right),$$

*the algorithm CH-PRFA returns a PFA realizing $\hat{p}$ such that*

$$\sum_{u \in \Sigma^t} |\hat{p}(u) - p(u)| \le \epsilon.$$

Before proving Proposition 6, we define few terms and prove a useful Lemma. We define,

$$\rho_0 = \|\hat{\boldsymbol{\alpha}}_0 - \boldsymbol{\alpha}_0\|_1,$$
$$\rho_\infty = \|\hat{\boldsymbol{\alpha}}_\infty - \boldsymbol{\alpha}_\infty\|_\infty,$$
$$\rho_\Sigma = \sum_{o \in \Sigma} \left\|\hat{A}_o - A_o\right\|_\infty = \left\|\hat{A} - A\right\|_\infty.$$

In addition, we introduce the following variables,

$$\gamma_k = \sum_{u \in \Sigma^k} \left\|\boldsymbol{\alpha}_0^\top A_u\right\|_1,$$
$$\gamma_\infty = \|\boldsymbol{\alpha}_\infty\|_\infty,$$
$$\gamma_\Sigma = \|A\|_\infty.$$

**Lemma 3.** *For a PFA (and so for a PRFA) $(\boldsymbol{\alpha}_0, A, \boldsymbol{\alpha}_\infty)$, we have*

$$\gamma_k \leq 1, \qquad\qquad \gamma_\infty \leq 1, \qquad\qquad \gamma_\Sigma \leq 1.$$

*Proof.*

$$\gamma_k = \sum_{u \in \Sigma^k} p(u\Sigma^\star) = p(\Sigma^k \Sigma^\star) \leq 1,$$
$$\gamma_\infty \leq \max_{u \in \mathcal{R}} \frac{p(u)}{p(u\Sigma^\star)} \leq 1,$$
$$\gamma_\Sigma = \left\|\sum_{o \in \Sigma} A_o \mathbf{1}\right\|_\infty = \|\mathbf{1} - \boldsymbol{\alpha}_\infty\|_\infty \leq 1.$$

$\square$

*Proof of Proposition 6.* Assuming that $\epsilon^{\text{est}} \leq \frac{\sigma_d^3}{648\sqrt{d}}$ we have by Proposition 2,

$$\epsilon^{\text{conv}} = \max_{u \in \hat{\mathcal{R}}} \left\|\hat{\mathbf{d}}_u - \mathbf{d}_u\right\|_2 \leq 162 \frac{\epsilon^{\text{est}}}{\sigma_d^2}.$$

Assuming that $\epsilon^{\text{conv}} \leq \frac{1}{\sqrt{d}}\kappa\sigma_d$, taking $\kappa = \frac{1}{3}$, setting $c_0, c_\infty$ and $c_\Sigma$ to be adequate constants and using that $d \geq 1$ and $|\Sigma| \geq 1$, we have by Propositions 3 to 5,

$$\rho_0 \leq \frac{9}{2}(\sqrt{2}+1)\frac{d^2}{\sigma_d^3}\epsilon^{\text{conv}} + 3\frac{\sqrt{d}}{\sigma_d}\epsilon^{\text{est}} \leq 9^3(\sqrt{2}+2)\frac{d^2}{\sigma_d^5}\epsilon^{\text{est}} + 3\frac{\sqrt{d}}{\sigma_d}\epsilon^{\text{est}} \leq \frac{c_0 d^2}{\sigma_d^5}\epsilon^{\text{est}}.$$

Likewise, if $\epsilon^{\text{conv}} \leq \frac{1}{\sqrt{d}}\kappa\sigma_d$ then,

$$\rho_\infty \leq \epsilon^{\text{conv}} \leq \frac{162}{\sigma_d^2}\epsilon^{\text{est}} = \frac{c_\infty}{\sigma_d^2}\epsilon^{\text{est}}.$$

Likewise, if $\epsilon^{\text{conv}} \leq \frac{1}{\sqrt{d}}\kappa\sigma_d$ then,

$$\rho_\Sigma \leq \frac{9}{2}(\sqrt{2}+2)\frac{d^2\sqrt{|\Sigma|}}{\sigma_d^3}\epsilon^{\text{conv}} + \frac{3}{2}\frac{d}{\sigma_d}\|\hat{\boldsymbol{\alpha}}_\infty - \boldsymbol{\alpha}_\infty\|_\infty + 3\frac{\sqrt{d}}{\sigma_d}\epsilon^{\text{conv}}$$
$$\leq 9^3(\sqrt{2}+2)\frac{d^2\sqrt{|\Sigma|}}{\sigma_d^5}\epsilon^{\text{est}} + 3^5\frac{d}{\sigma_d^3}\epsilon^{\text{est}} + 2 \cdot 3^5\frac{\sqrt{d}}{\sigma_d^3}\epsilon^{\text{est}}$$
$$\leq \frac{c_\Sigma d^2\sqrt{|\Sigma|}}{\sigma_d^5}\epsilon^{\text{est}}.$$

We denote $c = \max(c_0, c_\infty, c_\Sigma)$ and

$$\rho = \frac{cd^2 \sqrt{|\Sigma|}}{\sigma_d^5 (1-\kappa)^2} \epsilon^{\text{est}},$$

By the previous inequalities, we obtain

$$\max(\rho_0, \rho_\infty, \rho_\Sigma) \leq \rho.$$

Next, we can apply Lemma 5.4.4 in (Balle, 2013) that shows that for all integer $t \geq 0$,

$$\sum_{u \in \Sigma^t} |\hat{p}(u) - p(u)| \leq (\gamma_\infty + \rho_\infty) \left( (\gamma_\Sigma + \rho_\Sigma)^t \rho_0 + \rho_\Sigma \sum_{i=0}^{t-1} (\gamma_\Sigma + \rho_\Sigma)^i \gamma_{t-i-1} \right) + \gamma_t \rho_\infty.$$

By Lemma 3, we have

$$\sum_{u \in \Sigma^t} |\hat{p}(u) - p(u)| \leq (1 + \rho_\infty) \left( (1 + \rho_\Sigma)^t \rho_0 + \rho_\Sigma \sum_{i=0}^{t-1} (1 + \rho_\Sigma)^i \right) + \rho_\infty$$

$$= (1 + \rho_\infty) \left( 1 + (1 + \rho_\Sigma)^t \rho_0 + \rho_\Sigma \sum_{i=0}^{t-1} (1 + \rho_\Sigma)^i \right) - 1$$

$$= (1 + \rho_\infty) \left( 1 + (1 + \rho_\Sigma)^t \rho_0 + (1 + \rho_\Sigma)^t - 1 \right) - 1$$

$$= (1 + \rho_\infty)(1 + \rho_\Sigma)^t (1 + \rho_0) - 1.$$

Replacing $\rho_0$, $\rho_\Sigma$ and $\rho_\infty$ with $\rho$ leads to

$$\sum_{u \in \Sigma^t} |\hat{p}(u) - p(u)| \leq (1 + \rho)^{t+2} - 1.$$

Now we remark that if $p = \mathcal{O}(\frac{1}{t})$ then we can prove a bound on $\sum_{u \in \Sigma^t} |\hat{p}(u) - p(u)|$ without an exponential dependency on $t$ because $(1 + \frac{x}{2t})^t \leq 1 + x$ for $x \leq 1$. So, if

$$\epsilon^{\text{est}} \leq \frac{\sigma_d^3}{648\sqrt{d}}, \qquad \epsilon^{\text{conv}} \leq 3\frac{\sigma_d}{\sqrt{d}}, \qquad \rho \leq \frac{1}{2(t+2)},$$

then

$$\sum_{u \in \Sigma^t} |\hat{p}(u) - p(u)| \leq \left( 1 + \frac{2(t+2)\rho}{2(t+2)} \right)^{t+2} - 1 \leq 2(t+2)\rho.$$

As we have

$$\epsilon^{\text{conv}} \leq 162\frac{\epsilon^{\text{est}}}{\sigma_d^2}, \qquad \epsilon^{\text{est}} = \frac{\sigma_d^5}{cd^2\sqrt{|\Sigma|}}\rho,$$

the conditions are satisfied for

$$\epsilon^{\text{est}} \leq \frac{c'\sigma_d^5}{d^2\sqrt{|\Sigma|}(t+2)} \leq \min\left( \frac{1}{648} \frac{\sigma_d^3}{\sqrt{d}}, \frac{1}{2 \cdot 3^5} \frac{\sigma_d^3}{\sqrt{d}}, \frac{2}{81} \frac{\sigma_d^5}{cd^2\sqrt{|\Sigma|}(t+2)} \right), \tag{7}$$

where $c'$ is such that the last inequality is verified.

Finally, by proposition 1, with probability $1 - \delta$, we have for $|\mathcal{P}| \geq 2$ and for $n \geq 1$ that

$$\epsilon^{\text{est}} \leq \frac{1}{\sqrt{n}} \left( 1 + \sqrt{\log\left(\frac{|\mathcal{P}|}{\delta}\right)} \right) \leq \left( 1 + \sqrt{\frac{3}{2}} \right) \sqrt{\frac{1}{n} \log\left(\frac{|\mathcal{P}|}{\delta}\right)},$$

because $\log(2) > \frac{2}{3}$ and $\sqrt{\frac{3}{2}}\sqrt{\log\left(\frac{|\mathcal{P}|}{\delta}\right)} \geq 1$.

So, we can find a suitable constant $K$ such that $\epsilon$ if

$$n \geq K \frac{t^2 d^4 |\Sigma|}{\epsilon^2 \sigma_d^{10}} \log\left(\frac{|\mathcal{P}|}{\delta}\right),$$

then (7) holds and

$$\sum_{u \in \Sigma^t} |\hat{p}(u) - p(u)| \leq \frac{9}{4}(t+2)(t+1) \frac{cd^2 \sqrt{|\Sigma|}}{\sigma_d^5} \epsilon^{\text{est}} \leq \epsilon.$$

$\square$

*Proof of Theorem 1.* By Proposition 6, if we replace $\epsilon$ by $\frac{\epsilon}{t}$, there exists a suitable constant $K$ such that for all $t > 0$, $\epsilon > 0$, with probability $1 - \delta$, if

$$n \geq K \frac{t^4 d^4 |\Sigma|}{\epsilon^2 \sigma_d^{10}} \log\left(\frac{|\mathcal{P}|}{\delta}\right),$$

then

$$\sum_{u \in \Sigma^{\leq t}} |\hat{p}(u) - p(u)| \leq \frac{\epsilon}{t}.$$

Finally, we sum over $t$ to get

$$\sum_{u \in \Sigma^{\leq t}} |\hat{p}(u) - p(u)| \leq \epsilon.$$

$\square$

Here, we could continue, as it is done in (Bailly, 2011), to get a bound on

$$\sum_{u \in \Sigma^\star} |\hat{p}(u) - p(u)| \leq \sum_{u \in \Sigma^{\leq t}} |\hat{p}(u) - p(u)| + \hat{p}(\Sigma^{>t}) + p(\Sigma^{>t}),$$

by using the exponential decay of $\sum_{u \in \Sigma^{>t}} \hat{p}(u)$ and $\sum_{u \in \Sigma^{>t}} \hat{p}(u)$. Such a bound would depend on the spectral radius $\rho$ and $\hat{\rho}$ of $p$ and $\hat{p}$. Finally, this would give a bound in total variation. However, this technique does not give an explicit control over the constants involved.

## References

Bailly, Raphael. *Méthodes spectrales pour l'inférence grammaticale probabiliste de langages stochastiques rationnels.* PhD thesis, Université Aix Marseille 1, 2011.

Balle, Borja. *Learning finite-state machines: statistical and algorithmic aspects.* PhD thesis, Universitat Politècnica de Catalunya, 2013.

Denis, François, Gybels, Mattias, and Habrard, Amaury. Dimension-free concentration bounds on Hankel matrices for spectral learning. In *Proceedings of the 31$^{th}$ International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, pp. 449–457, 2014.

Gillis, N. and Vavasis, S. A. Fast and robust recursive algorithmsfor separable nonnegative matrix factorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(4):698–714, April 2014. ISSN 0162-8828.

Hsu, Daniel, Kakade, Sham M, and Zhang, Tong. A spectral algorithm for learning hidden Markov models. *Journal of Computer and System Sciences*, 78(5):1460–1480, 2012.

Lőtstedt, Per. Perturbation bounds for the linear least squares problem subject to linear inequality constraints. *BIT Numerical Mathematics*, 23(4):500–519, 1983. ISSN 0006-3835. doi: 10.1007/BF01933623.

McDiarmid, Colin. On the method of bounded differences. *Surveys in Combinatorics*, 141(1):148–188, 1989.

Wedin, Per-Åke. Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics*, 12 (1):99–111, 1972.