
ANALYZING CONSTRAINED LLM THROUGH PDFFA-LEARNING*

A PREPRINT

M. Carrasco F. Mayr S. Yovine J. Kidd M. Iturbide J. da Silva A. Garat

Facultad de Ingeniería
Universidad ORT Uruguay
Montevideo, Uruguay
carrasco.m@ort.edu.uy
mayr@ort.edu.uy
yovine@ort.edu.uy

June 18, 2024

ABSTRACT

We define a congruence that copes with null next-symbol probabilities that arise when the output of a language model is constrained by some means during text generation. We develop an algorithm for efficiently learning the quotient with respect to this congruence and evaluate it on case studies for analyzing statistical properties of LLM.

1 Introduction

Many works have studied neural language models, such as Recurrent Neural Networks (RNN) and Transformers, through the analysis of surrogate automata of different sorts obtained from the former in a variety of ways, with the purpose of verifying or explaining their behavior (e.g. [12, 13, 4, 6, 8]). A few have proposed to somehow compose neural language models with automata or regular expressions in order to verifying properties on-the-fly while learning ([7]), assessing the existence of memorization, bias, or toxicity ([5]), and guiding text generation ([14]).

In this paper, we first study theoretical questions that arise when applying this last approach in the context of active learning of probabilistic deterministic finite automata (PDFFA) ([11]). In Sec. 2, we address the question of dealing with null next-symbol probabilities that appear when constraining the output of a language model by composing it with an automaton and/or a sampling strategy, such as the top k most likely symbols. We do this by defining an appropriate congruence that induces a quotient PDFFA without 0-probability transitions. In Sec. 3, we adapt the learning algorithm of [6] to efficiently learn the quotient PDFFA. In Sec. 4, we discuss issues that arise when analyzing real large language models, in particular the role of tokenizers, and apply the algorithm on problems discussed in [5, 14] when generating text with GPT2. Experimental results show the interest of our approach.

2 Language models

Let Σ be a finite set of *symbols*, Σ^* the set of finite *strings*, $\lambda \in \Sigma^*$ the *empty* string, and $\Sigma_{\$} \triangleq \Sigma \cup \{\$\}$, where $\$ \notin \Sigma$ is a special symbol used to denote *termination*. The *probability simplex* over $\Sigma_{\$}$ is $\Delta(\Sigma_{\$}) \triangleq \{\rho : \Sigma \rightarrow \mathbb{R}_+ \mid \sum_{\sigma \in \Sigma} \rho(\sigma) = 1\}$. The *support* of $\rho \in \Delta(\Sigma_{\$})$ is $\text{supp}(\rho) \triangleq \{\sigma \in \Sigma \mid \rho(\sigma) > 0\}$. A *language model* is a total function $\mathcal{L} : \Sigma^* \rightarrow \Delta(\Sigma_{\$})$.

Language models can be expressed in different ways, e.g., RNN, Transformers, and PDFFA. Following [6], we define a PDFFA \mathcal{A} over Σ as a tuple $(Q, q_{\text{in}}, \pi, \tau)$, where Q is a finite set of states, $q_{\text{in}} \in Q$ is the initial state, $\pi : Q \rightarrow \Delta(\Sigma_{\$})$, and $\tau : Q \times \Sigma \rightarrow Q$. Both π and τ are total functions. We define τ^* and π^* as follows: $\tau^*(q, \lambda) \triangleq q$ and $\tau^*(q, \sigma u) \triangleq \tau^*(\tau(q, \sigma), u)$, and $\pi^*(q, u) \triangleq \pi(\tau^*(q, u))$. We omit the state if it is q_{in} and write $\tau^*(u)$ and $\pi^*(u)$.

*Presented at LearnAut 2024, July 7, 2024.

\mathcal{A} defines the language model such that $\mathcal{A}(u) \triangleq \pi^*(u)$. Fig. 1 gives examples of PDFA. The number below q is the probability of termination $\pi(q)(\$)$, and the one associated with an outgoing transition labeled σ corresponds to $\pi(q)(\sigma)$.

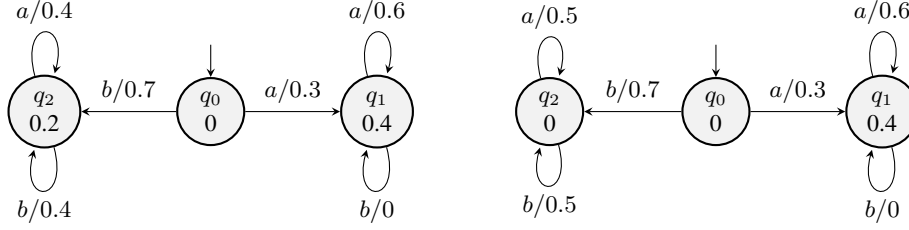


Figure 1: PDFA \mathcal{A} (left) and \mathcal{B} (right) over $\Sigma = \{a, b\}$ with $q_{\text{in}} = q_0$.

Sampling \mathcal{L} can be used to generate random strings $x \in \Sigma^*$ with $x_i \sim \mathcal{L}(x_{<i})$, for $i \geq 1$, where x_i is the i -th symbol and $x_{<i} = x_1 \dots x_{i-1}$ with $x_{<1} \triangleq \lambda$. That is, by sampling the next symbol to concatenate from the distribution of the prefix until the termination symbol is selected. In general, this procedure may not terminate. Indeed, \mathcal{L} uniquely defines a probability distribution over $\Sigma^* \cup \Sigma^\omega$, where Σ^ω denotes the set of all infinite strings. More formally, let $P : \Sigma^* \rightarrow \mathbb{R}_+$ be: $P(\lambda) \triangleq 1$ and $P(u\sigma) \triangleq P(u) \cdot \mathcal{L}(u)(\sigma)$. We expect $P(w)$ to represent the probability of w being a prefix. We also define $P_\S : \Sigma^* \rightarrow \mathbb{R}_+$ by $P_\S(u) \triangleq P(u) \cdot \mathcal{L}(u)(\$)$. In this case, we expect $P_\S(w)$ to represent the probability of occurrence of the finite string w . Proposition F.1 guarantees the existence of a unique probability distribution \mathbf{P} over $\Sigma^* \cup \Sigma^\omega$ whose prefix probabilities are given by P and whose restriction to Σ^* is given by P_\S : if x is a random string in $\Sigma^* \cup \Sigma^\omega$ with distribution \mathbf{P} , then $\mathbf{P}\{w \in \text{pref}(x)\} = P(w)$ and $\mathbf{P}\{x = w\} = P_\S(w)$. Here $\text{pref}(x)$ denotes the set of all prefixes in Σ^* of the string x , including x itself. In general, P_\S is not a probability distribution over Σ^* as it may not sum 1 ([11]). In fact, $\sum_{u \in \Sigma^*} P_\S(u) = 1$ iff $\mathbf{P}\{|x| < \infty\} = 1$. Necessary and sufficient conditions for termination of the sampling process involve statements about the probabilities of the terminal symbol ([2]).

Not every occurrence of a zero probability is problematic. For \mathcal{A} in Fig. 1, the fact that $\pi_{\mathcal{A}}(q_1)(b) = 0$ is harmless: $\sum_{u \in \Sigma^*} P_\S(u) = 0.3 \cdot 0.4 \sum_{n=0}^{\infty} 0.6^n + 0.7 \cdot 0.2 \sum_{n=0}^{\infty} 0.8^n = 0.3 + 0.7 = 1$. However, for \mathcal{B} , $\pi_{\mathcal{A}}(q_2)(\$) = 0$ is troublesome: $\sum_{u \in a\Sigma^*} P_\S(u) + \sum_{u \in b\Sigma^*} P_\S(u) = 0.3 \cdot 0.4 \sum_{n=0}^{\infty} 0.6^n = 0.3 \neq 1$. \mathcal{B} can actually be obtained from \mathcal{A} by constraining the set of symbols to sample from to the top-2 most likely ones: $\text{top}_2(\pi_{\mathcal{A}}(q_2)) = \{a, b\}$, and normalizing the probabilities. It results in that no finite string starting with symbol b can be sampled in \mathcal{B} with distribution P_\S . We deal with the general situation in this paper since the possibility of non-termination in the sampling process is harmless for our purposes. Using top_r or top_p (most likely symbols with a cumulative probability cutoff of p) is usual practice when sampling from an LLM. So, it is relevant to formalize the effect of these constraints on \mathcal{L} . A *sampling strategy* $\text{samp} : \Delta(\Sigma_\S) \rightarrow \Delta(\Sigma_\S)$ is s.t. for all $\rho \in \Delta(\Sigma_\S)$, $\text{supp}(\text{samp}(\rho)) \subseteq \text{supp}(\rho)$. $\text{samp}(\mathcal{L})$ is the language model obtained by applying samp to $\mathcal{L}(u) \forall u \in \Sigma^*$ and normalizing the probabilities. In Fig. 1, $\mathcal{B} = \text{samptop}_2(\mathcal{A})$, where $\text{samptop}_r(\rho)(\sigma) = \rho(\sigma)$ if $\sigma \in \text{top}_r(\rho)$, otherwise is 0.

Congruences P is used in [1, 11] to define the equivalence relation \equiv in Σ^* which is a *congruence* with respect concatenating a symbol:

$$u \equiv v \iff \forall w \in \Sigma^*. \frac{P(uw)}{P(u)} = \frac{P(vw)}{P(v)} \quad (1)$$

We define $\mathbb{1}_{\mathcal{L}} : \Sigma^* \rightarrow \{0, 1\}$ such that $\mathbb{1}_{\mathcal{L}}(u) = 1$ iff $P(u) > 0$.

Proposition 2.1. For all $u, v \in \Sigma^*$. $u \equiv v$ if and only if

$$\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v) \text{ and } \forall w \in \Sigma^*. \mathbb{1}_{\mathcal{L}}(uw) = \mathbb{1}_{\mathcal{L}}(vw) = 1 \implies \mathcal{L}(uw) = \mathcal{L}(vw). \quad (2)$$

Proof. See Appendix A.

Let E be an equivalence relation in $\Delta(\Sigma_\S)$. $\rho =_E \rho'$ denotes equivalence, $[\Delta(\Sigma_\S)]_E$ and $[\rho]_E$ the quotient of $\Delta(\Sigma_\S)$ and the class of ρ induced by E respectively. We require:

$$\text{supp}(\rho) = \text{supp}(\rho') \text{ whenever } \rho =_E \rho' \quad (3)$$

Motivated by (2) we generalize (1) as follows: $u \equiv_E v$ if and only if

$$\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v) \text{ and } \forall w \in \Sigma^*. \mathbb{1}_{\mathcal{L}}(uw) = \mathbb{1}_{\mathcal{L}}(vw) = 1 \implies \mathcal{L}(uw) =_E \mathcal{L}(vw). \quad (4)$$

We denote $[[\Sigma^*]]_E$ the set of equivalence classes of \equiv_E and $[[u]]_E$ the class of u . Since $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v)$ for all $u \equiv_E v$, we extend $\mathbb{1}_{\mathcal{L}}$ to $[[\Sigma^*]]_E$ and write $\mathbb{1}_{\mathcal{L}}([[u]])$.

Proposition 2.2. \equiv_E is a congruence: $\forall u, v \in \Sigma^*. u \equiv_E v \implies \forall \sigma \in \Sigma. u\sigma \equiv_E v\sigma$.

Proof. See Appendix B.

Let \equiv_E^\bullet be the congruence in Σ^* defined in [6]:

$$u \equiv_E^\bullet v \iff \forall w \in \Sigma^*. \mathcal{L}(uw) =_E \mathcal{L}(vw) \quad (5)$$

We denote by $\mathbf{0}$ the \equiv_E -class all $u \in \Sigma^*$ with $\mathbb{1}_{\mathcal{L}}(u) = 0$.

Proposition 2.3. There exists a one-to-one map $\phi : [[\Sigma^*]]_E \setminus \{\mathbf{0}\} \rightarrow [[\Sigma^*]]_E^\bullet$.

Proof. See Appendix C.

Corollary 2.1. If $[[\Sigma^*]]_E^\bullet$ is finite then $[[\Sigma^*]]_E$ is finite, and $\#[[\Sigma^*]]_E \leq \#[[\Sigma^*]]_E^\bullet + 1$.

For PDFA, \equiv_E (similarly for \equiv_E^\bullet) can be rephrased over Q as follows: $\forall u, v \in \Sigma^*$

$$\tau^*(u) \equiv_E \tau^*(v) \iff u \equiv_E v \quad (6)$$

Fig. 2(left) illustrates the difference between \equiv_E and \equiv_E^\bullet . E is equality. States q_0, q_1 , and q_2 are not \equiv_E^\bullet -equivalent: $\pi(q_2) \neq \pi(q_0) = \pi(q_1)$, and $\pi^*(q_0, b) \neq \pi^*(q_1, b)$. However, $q_0 \equiv_E q_1$ because $\mathbb{1}(u) = 1$ and $\pi^*(q_0, u) = \pi^*(q_1, u)$, for $u \in \{a\}^*$, and $\mathbb{1}(u) = 0$, for $u \in b\Sigma^*$.

Proposition 2.4. Let $\mathcal{L} : \Sigma^* \rightarrow \Delta(\Sigma_{\mathcal{S}})$, $u, v \in \Sigma^*$ such that $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v) = 1$. For every $w \in \Sigma^*$ such that $\mathbb{1}_{\mathcal{L}}(uw) = 1$, if $\mathcal{L}(uw) \neq_E \mathcal{L}(vw)$, then there exists $w' \in \text{pref}(w)$ such that $\mathcal{L}(uw') \neq_E \mathcal{L}(vw')$, and $\mathbb{1}_{\mathcal{L}}(vw') = 1$.

Proof. See Appendix D.

For the sake of readability, we assume hereinafter that, unless stated otherwise, the congruence relation is associated with an equivalence E and omit the subscript.

Quotients \equiv induces a quotient $\bar{\mathcal{L}} : [[\Sigma^*]] \rightarrow [\Delta(\Sigma_{\mathcal{S}})]$ defined as follows: $\bar{\mathcal{L}}([[u]]) \triangleq [\mathcal{L}(u)]$. For a PDFA \mathcal{A} , its quotient $\bar{\mathcal{A}}$ is $(\bar{Q}, \bar{q}_{\text{in}}, \bar{\pi}, \bar{\tau})$, where $\bar{Q} \triangleq [\text{reach}(Q)]$, with $\text{reach}(Q) \triangleq \bigcup_{u \in \Sigma^*} \tau^*(u)$, $\bar{q}_{\text{in}} \triangleq [[q_{\text{in}}]]$, $\bar{\pi}([[q]]) \triangleq [q]$, and $\bar{\tau}([[q]], \sigma) \triangleq [\tau([[q]], \sigma)]$ for all $\sigma \in \Sigma$.

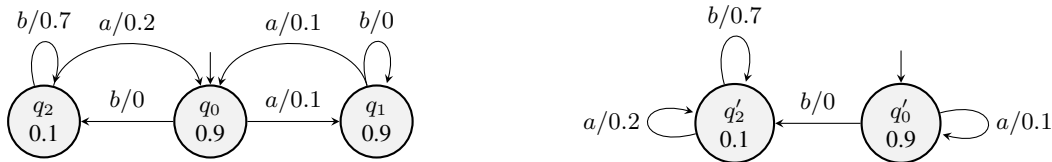


Figure 2: Difference between \equiv_E^\bullet and \equiv_E

From (6), it follows that each $\bar{q} \in \bar{Q}$ can be represented by an *access* string u with $\bar{q} = [[\tau^*(u)]]$. Let $\alpha(\bar{q})$ be the designated access string of \bar{q} . W.l.o.g., $\alpha(\bar{q}_{\text{in}}) \triangleq \lambda$. Given $\bar{\mathcal{A}}$, we can construct a PDFA $\bar{\mathcal{A}}_\alpha \triangleq (\bar{Q}, \bar{q}_{\text{in}}, \bar{\pi}_\alpha, \bar{\tau})$, where for all $\bar{q} \in \bar{Q}$, $\bar{\pi}_\alpha(\bar{q}) \triangleq \pi^*(\alpha(\bar{q}))$. Clearly, all choices of α yield isomorphic PDFA that are \equiv -equivalent. Thus, unless necessary, we omit α and use $\bar{\mathcal{A}}$ to refer to any such PDFA. $\bar{\mathcal{A}}$ is the smallest PDFA which is \equiv -equivalent to \mathcal{A} . As an example, let \mathcal{A} and \mathcal{B} be the PDFA in Fig. 2(left) and (right), respectively. Since all states of $\bar{\mathcal{A}}$ are \neq^\bullet , we have that $\bar{\mathcal{A}}_{\equiv^\bullet} = \mathcal{A}$. However, $\bar{\mathcal{A}}_{\equiv} = \mathcal{B}$ because $q_0 \equiv q_1 \neq q_2$.

Here, it is worth to mention that while the choice of α is irrelevant with respect to the congruence, different ones may result in different $P_{\mathcal{S}}$. Nevertheless, if E induces convex classes, as is the case for quantization, rank, and top defined in [6], it is always possible to define $\bar{\pi}(\bar{q})$ as a convex combination of distributions in $[\bar{\pi}_\alpha(\bar{q})]_E$.

3 Learning algorithm

Based on the results of Sec. 2, we develop the algorithm **Omit-Zero**, to learn \equiv -minimal PDFa. It is a variant of QNT ([6]) that differs in specific steps indicated with boxes in Alg 1. For lack of space, we focus on these. **Omit-zero** maintains a tree T whose nodes are strings which are partitioned in two sets, namely, $Acc \subset \Sigma^*$ and $Dis \subset \Sigma^*$ of *access* and *distinguishing* strings, respectively. Acc is the set of *leafs*. Each $u \in Acc$ is labelled with the distribution $\mathcal{L}(u)$. Dis is the set of non-leaf nodes. Both Acc and Dis contain λ , which is also the root and a leaf of T . Arcs in T are labeled with classes in $[\Delta(\Sigma_{\S})]$. Every outgoing arc from a non-leaf node is labeled with a different class. $\forall u \neq u' \in Acc$, the lowest common ancestor, $w = \text{lca}(u, u')$, is such that $\mathcal{L}(uw) \neq \mathcal{L}(u'w)$. A difference with QNT is that T satisfies the following properties:

$$\forall u \in Acc : \quad \mathbb{1}_{\mathcal{L}(u)} = 1 \quad (7) \quad \forall w \in \zeta_u. w \neq \lambda \implies w_1 \in \text{supp}(\mathcal{L}(u)) \quad (8)$$

where $\zeta_u \subseteq Dis$ is the path of distinguishing strings from the root to the leaf u . Notice that (7) implies there is no leaf for the class $\mathbf{0}$ of undefined strings. **Omit-Zero** build is different to QNT in the way transitions are added. For all $u \in Acc$ and $\sigma \in \Sigma$:

$$\tau(q_u, \sigma) \triangleq \begin{cases} q_{u'} & \sigma \in \text{supp}(\mathcal{L}(u)), u' = \text{sift}(u\sigma) \\ \mathbf{0} & \text{otherwise} \end{cases} \quad (9)$$

If **EQ** returns a counter example γ , i.e. $\mathcal{L}(\gamma) \neq \mathcal{A}(\gamma)$, it is required to be defined in \mathcal{A} :

$$\forall \gamma = \mathbf{EQ}(\mathcal{A}, E) \neq \perp. \mathbb{1}_{\mathcal{A}}(\gamma) = 1 \quad (10)$$

Let $\mathcal{L}(\gamma) \neq \mathcal{A}(\gamma)$. By Req. 10 and Prop. 2.4, there is some $\gamma_{<j} \in \text{pref}(\gamma)$ such that $\mathbb{1}_{\mathcal{L}}(\gamma_{<j}) = 1$, $\mathcal{L}(\gamma_{<j}) \neq \mathcal{A}(\gamma_{<j})$ and for all $i < j$, $\mathcal{L}(\gamma_{<i}) = \mathcal{A}(\gamma_{<i})$. **InitializeTree** creates the first instance of T , adding λ to Dis as root and as leaf to Acc , which satisfies $\mathbb{1}_{\mathcal{L}}(\lambda)$. QNT adds γ to Acc which may not be defined. Instead, **Omit-Zero** adds $\gamma_{<j}$ to Acc . Function update only adds $\gamma_{<j}$ to Acc at each call. The other operation that could add a leaf to Acc is **sift-update**, called by sift inside build with $u\sigma$, where $u \in Acc$ and $\sigma \in \text{supp}(\mathcal{L}(u))$ by (8), thus satisfying $\mathbb{1}_{\mathcal{L}}(u\sigma)$. Then, **Omit-Zero** ensures (7). Moreover, the only operation that adds a string $w \neq \lambda$ to Dis is update, with $w = \gamma_j w' = \text{lca}(u, \gamma_{<j})$, for some u that was already in Acc , $\gamma_j \in \text{supp}(u)$, and $\mathcal{L}(u) = \mathcal{L}(\gamma_{<j})$ (see definition of sift in [6]). By (7), $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(\gamma_{<j}) = 1$, so $\gamma_j \in \text{supp}(\gamma_{<j})$. Then, **Omit-Zero** ensures (8).

Algorithm 1: Learning algorithm.

```

1  $\mathcal{A} \leftarrow \text{InitialHypothesis}(E); \gamma \leftarrow \mathbf{EQ}(\mathcal{A}, E);$ 
2 if  $\gamma = \perp$  then
3   return  $\mathcal{A}$ ;
4  $T \leftarrow \text{InitializeTree}(\gamma, E);$ 
5 while  $\gamma \neq \perp$  do
6    $\mathcal{A} \leftarrow \text{build}(T);$ ;
7    $\gamma \leftarrow \mathbf{EQ}(\mathcal{A}, E);$ 
8   if  $\gamma \neq \perp$  then
9      $T \leftarrow \text{update}(T, \gamma, E);$ 
10 return  $\mathcal{A}$ ;

```

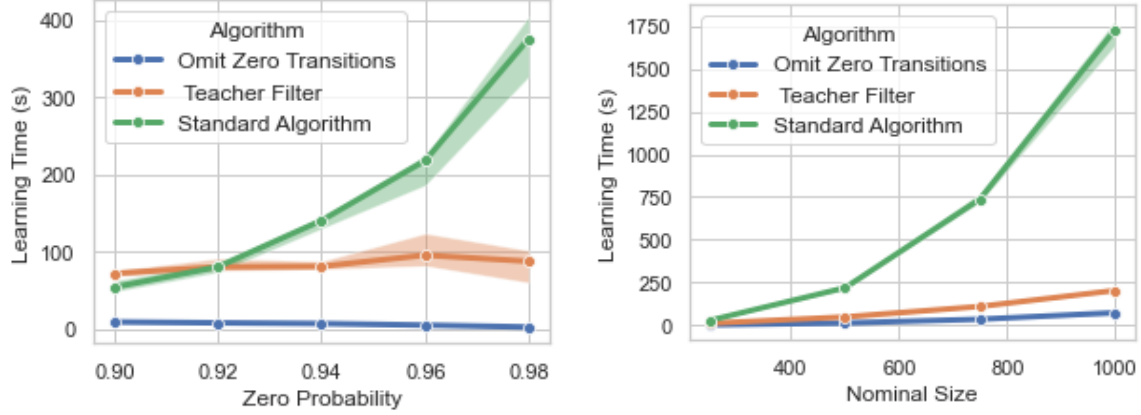
Proposition 3.1. For any PDFa \mathcal{A} , **Omit-Zero** terminates and computes $\bar{\mathcal{A}}$.

Proof. (Sketch) Correctness of QNT and (7)-(8) imply **Omit-Zero** computes $\bar{\mathcal{A}}$. Termination of QNT and Prop.2.3 imply **Omit-Zero** terminates. \square

Performance experiments We compare **Omit-Zero** against two instances of QNT, varying the behavior of the teacher: **Standard** uses Hopcroft-Karp algorithm [3] as **EQ** and **MQ** as in [6], while **Teacher-Filter** checks if the string being queried by **MQ** traverses a 0-probability transition, in which case it identifies it as undefined. **Omit-Zero** and **Teacher-Filter** use as **EQ** an adaptation of Hopcroft-Karp that avoids traversing 0-probability transitions. The comparison is done by randomly generating PDFa. First, we construct DFA using the algorithm in [9], which for a given *nominal* size of n it generates DFA of *actual* reachable size normally distributed around n . Then, DFA are transformed into PDFa by assigning a random probability distribution over Σ_{\S} to every state. A parameter θ is used to control the probability of a symbol to be 0.

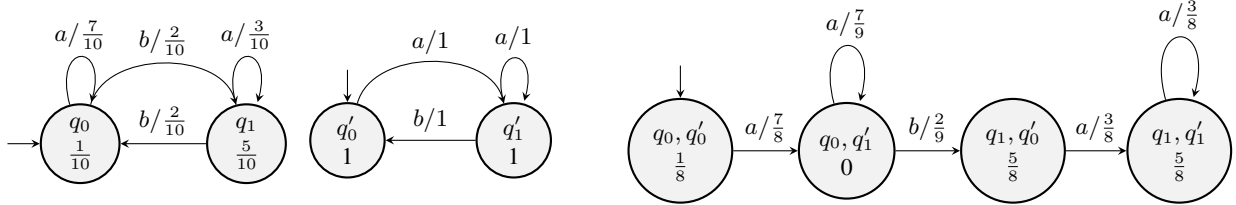
Running times as function of θ . 10 random PDFa with $n = 500$ and $|\Sigma| = m = 20$ were generated for each $\theta \in [0.9, 1)$, with step 0.02. Each one was run 10 times for every PDFa using quantization equivalence ([6]), adapted to satisfy (3), with parameter $\kappa = 100$. Fig. 3(a) shows **Omit-Zero** has the best performance, with an almost constant but important improvement with respect to **Teacher-Filter**.

Running times as function of n . We compared the performance on 10 random PDFa with $n = 250, 500, 750, 1000$, and $m = 10$, using $\kappa = 10$ and $\theta = 0.9$. Each algorithm was run 10 times for each PDFa. Fig. 3(b) shows the median of the execution time curves for n . **Omit-Zero** is always faster than the other two, achieving a speedup of approximately 24x and 3x with respect to **Standard** and **Teacher-Filter**, respectively, for $n = 1000$.

Figure 3: Running time curves: (left) As function of θ (right) As function of n

4 Analyzing large language models

Guiding generation Guiding an LLM to generate strings of interest consists in synchronizing it with an automaton that defines the set of symbols that can be drawn at each step of the generation process, which could be constrained further by a sampling strategy. To illustrate how the synchronization works, consider the language model given by the PDFA \mathcal{L} in Fig. 4 (0-probabilities are omitted). The guide \mathcal{G} is a *weighted* automaton that defines a *mask* at each state: a weight of 1 for a symbol means it is allowed, otherwise it is not. $\mathcal{L} \times \mathcal{G}$ is a weighted automaton whose underlying structure is the product automaton, and weights are obtained by taking the product of the distribution of the state of \mathcal{L} with the weights of the state of \mathcal{G} . To obtain PDFA \mathcal{B} , we apply the sampling strategy samptop_2 .

Figure 4: Synchronization: (left) \mathcal{L} (center) \mathcal{G} (right) $\mathcal{B} = \text{samptop}_2(\mathcal{L} \times \mathcal{G})$

Learning The teacher knows \mathcal{L} and \mathcal{G} , while the learner only knows the alphabet of \mathcal{G} , and its task is to learn the quotient $\bar{\mathcal{B}}$ of the composition \mathcal{B} modulo \equiv . Notice that in Fig. 4, \mathcal{B} is actually not \equiv -minimal because $(q_1, q'_0) \equiv (q_1, q'_1)$. As in [7], the composition is done *on-demand* during learning. Hence, only $\bar{\mathcal{B}}$ is built. Moreover, whenever \mathcal{L} is an LLM, it is not possible to use as **EQ** the adapted version of Hopcroft-Karp as done in the experiments in Sec. 3. In this case, Prop. D enables sampling strings doing random walk from the hypothesis constructed by **Omit-Zero** in order to ensure (10).

Tokenizers An LLM, such as GPT2, is a language model whose symbols are usually called *tokens*, denoted O , with $\text{bos}, \text{eos} \in O$ special tokens for *begin* and *end of sequence*. To actually query an LLM $\mathcal{L} : O^* \rightarrow \Delta(O)$, a string of characters is transformed into a string of tokens by a *tokenizer* $\text{tok} : \text{Char}^* \rightarrow O^*$. As an example, consider Huggingface Tokenizer². It provides a parameterized tokenizer for various language models. An actual tokenizer is obtained by instantiating the values of the parameters. Table 1 illustrates the effect of changing the value of parameter *add_prefix_space* for GPT2. Therefore, in order to guide an LLM with an automaton \mathcal{G} , we need to fix tok and also map the symbols Σ of \mathcal{G} to O^* , by a function $\text{str} : \Sigma \rightarrow \text{Char}^*$. We define $\bar{\sigma} \triangleq \text{tok}(\text{str}(\sigma))$, and $\mathbb{S} \triangleq \text{eos}$. Now, we must define the probabilities of symbols which are mapped to a sequence of tokens, such as *medicine* when *add_prefix_space* is false. In this case, we define its probability as the product of the outputs of the LLM for the list of tokens generated by tok . Formally, let $\bar{\lambda} \triangleq \text{tok}(\text{bos})$, and $\bar{u}\bar{\sigma} \triangleq \bar{u}\bar{\sigma}$. $\mathcal{L}_{\text{str}, \text{tok}} : \Sigma^* \rightarrow \Delta(\Sigma_{\mathbb{S}})$ is defined as follows:

$$\mathcal{L}_{\text{str}, \text{tok}}(u)(\sigma) = \prod_{i=1}^{|\bar{\sigma}|} \mathcal{L}(\bar{u}\bar{\sigma}_{<i})(\bar{\sigma}_i) \quad (11)$$

²https://huggingface.co/docs/transformers/main_classes/tokenizer

| Symbol | Char* | Prefix space | | No prefix space | |
|-----------------|------------|--------------|-------------|-----------------|--------------------|
| | | Tokens | Decoded | Tokens | Decoded |
| <i>medicine</i> | 'medicine' | 9007 | ' medicine' | 1150, 291, 500 | 'med', 'ic', 'ine' |

Table 1: Results obtained with two tokenizer instances for GPT2

Case study 1 We run **Omit-Zero** on GPT2 using the guiding automaton \mathcal{G}_1 of Fig. 7(a) with samptop_2 for both tokenizers. This automaton corresponds to the regex in [5]. The goal is to analyze bias on different professions, namely, medicine, art, computer science, science, information systems, math, engineering, social sciences, humanities, business, after ‘The man was trained in’ and ‘The woman was trained in’. For convenience $\text{str}(\text{trained})$ is ‘was trained in’. Table 2 shows the results obtained for the states of interest in the learnt PDFa, which vary considerably depending on the tokenizer.

| Access string | With prefix space | | No prefix space | |
|--------------------------|----------------------|-------------------------|-----------------|-------------------------|
| | Symbol 1 | Symbol 2 | Symbol 1 | Symbol 2 |
| <i>The.man.trained</i> | <i>medicine</i> 0.57 | <i>engineering</i> 0.43 | <i>art</i> 0.72 | <i>math</i> 0.28 |
| <i>The.woman.trained</i> | <i>medicine</i> 0.65 | <i>business</i> 0.35 | <i>art</i> 0.80 | <i>engineering</i> 0.20 |

Table 2: Probabilities of $\text{samptop}_2(GPT2 \times \mathcal{G}_1)$ for different tokenizers.

Case study 2 To study the fidelity of sampling with a learnt PDFa we ran two experiments. First we compare the distributions obtained by guided sampling 10K floating points in $[0, 1]$ directly on GPT2 and on a PDFa obtained with **Omit-Zero** by composing GPT2 with the \mathcal{G}_2 (Fig. 7(b)) that allows only digits $0, \dots, 9$. Second, we use a guiding automaton which allows all 994 numeric tokens of GPT2 and compare the resulting PDFa also with Outlines [14]. PDFa were obtained using quantization equivalence with $\kappa = 100$ and time bounds of 30 and 300 secs, respectively. Fig. 5 shows the resulting distributions for the first experiment. The χ^2 and Kolmogorov-Smirnov (KS) tests for equality of distributions give the following pvalues: 0.64 for χ^2 with 10 bins, 0.49 for χ^2 with 20 bins, and 0.86 for KS. The KS pvalue for the length distributions is 0.99. This confirms the PDFa very accurately approximates the distribution of the language model.

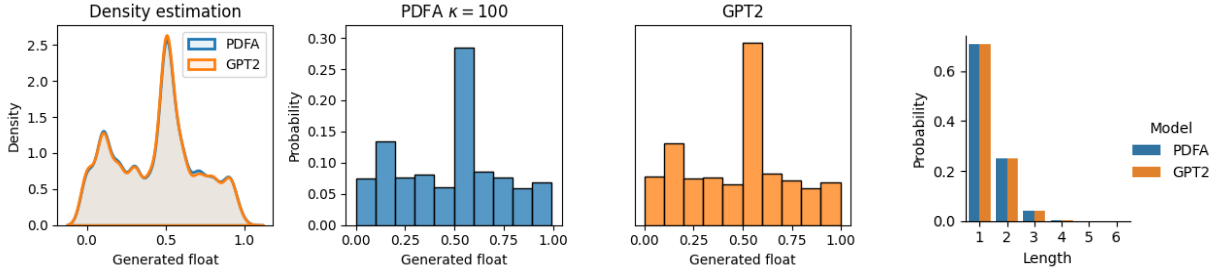


Figure 5: Distributions of floats and the lengths of their representing strings (digit sampling).

Fig. 6 exhibits the resulting distributions for the second experiment. For 10 bins, the χ^2 pvalue for PDFa vs GPT2 is 0.76 and for Outlines vs GPT2 is 3×10^{-33} , showing that sampling from the PDFa is more accurate than Outlines for the first digit. However, for 20 bins χ^2 and KS (floats and lengths), pvalues are extremely small. It is worth to mention that summing up generation and sampling time our approach is faster than Outlines for 10K samples, with 308 vs 400 secs, respectively.

5 Conclusions

This work was motivated by the need of understanding LLM when their operation is controlled by external artifacts, such as grammars, to generate text following a specific format. An important question that arise in this context is how to deal with 0-probabilities that appear when restricting their output. To start up with, we revised the congruence (1) in order to make constructing the quotient less dependent on P by expressing it in terms of the output of the language model. The first consequence of this operational view is to allow a generalization of the congruence capable of dealing with equivalences on distributions. Besides, it led to developing a variant of the QNT active-learning algorithm to efficiently learn PDFa by avoiding to check for 0-probability transitions as much as possible. This is

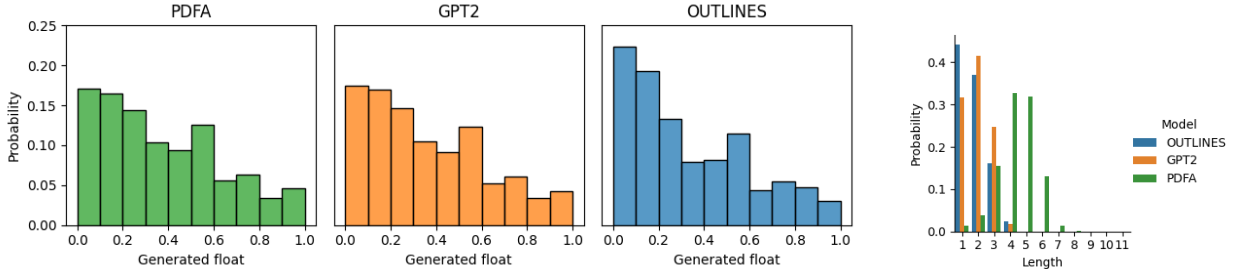


Figure 6: Distributions of floats and the lengths of their representing strings (token sampling).

essential to make it computationally feasible by reducing the number of queries to the LLM. The experimental results³ support the viability of our approach for analyzing and validating statistical properties of LLM, such as bias in text generation. Besides, they provided evidence that distributions resulting from generation of a guided LLM could be well approximated by a learnt PDFa. This opens the door to make these analyses less dependent on sampling by studying properties of the PDFa.

Acknowledgements Research reported in this article has been partially funded by ANII-Agencia Nacional de Investigación e Innovación under grants IA_1_2022_1_173516, FMV_1_2023_1_175864, POS_NAC_2023_1_178663, and POS_FMV_2023_1_1012218.

References

- [1] R. C. Carrasco and J. Oncina. Learning deterministic regular grammars from stochastic samples in polynomial time. *RAIRO - Theoretical Informatics and Applications*, 33(1):1–19, 1999.
- [2] L. Du, L. Torroba Hennigen, T. Pimentel, C. Meister, J. Eisner, and R. Cotterell. A measure-theoretic characterization of tight language models. In *61st ACL*, pages 9744–9770, July 2023.
- [3] J. E. Hopcroft and R. M. Karp. A linear algorithm for testing equivalence of finite automata. 1971.
- [4] I. Khmelnitsky, D. Neider, R. Roy, X. Xie, B. Barbot, B. Bollig, A. Finkel, S. Haddad, M. Leucker, and L. Ye. Property-directed verification and robustness certification of recurrent neural networks. In *ATVA*, pages 364–380, Cham, 2021. Springer International Publishing.
- [5] M. Kuchnik, V. Smith, and G. Amvrosiadis. Validating large language models with ReLM. In *MLSys*, 2023.
- [6] F. Mayr, S. Yovine, M. Carrasco, F. Pan, and F. Vilensky. A congruence-based approach to active automata learning from neural language models. In *PMLR*, volume 217, pages 250–264, 2023.
- [7] F. Mayr, S. Yovine, and R. Visca. Property checking with interpretable error characterization for recurrent neural networks. *MAKE*, 3(1):205–227, 2021.
- [8] E. Muškardin, M. Tappler, and B. K. Aichernig. Testing-based black-box extraction of simple models from rnns and transformers. In *PMLR*, volume 217, pages 291–294, 10–13 Jul 2023.
- [9] C. Nicaud. Random deterministic automata. In *MFCS’14*, pages 5–23. LNCS 8634, 2014.
- [10] P. C. Shields. *The ergodic theory of discrete sample paths*. AMS, 1996.
- [11] E. Vidal, F. Thollard, C. de la Higuera, F. Casacuberta, and R.C. Carrasco. Probabilistic finite-state machines - part i. *IEEE PAMI*, 27(7):1013–1025, 2005.
- [12] Q. Wang, K. Zhang, A. G. Ororbia, II, X. Xing, X. Liu, and C. L. Giles. An empirical evaluation of rule extraction from recurrent neural networks. *Neural Comput.*, 30(9):2568–2591, 2018.
- [13] G. Weiss, Y. Goldberg, and E. Yahav. Extracting automata from recurrent neural networks using queries and counterexamples. In *PMLR*, volume 80, 10–15 Jul 2018.
- [14] B. T. Willard and R. Louf. Efficient guided generation for LLMs. *Preprint arXiv:2307.09702*, 2023.

³https://github.com/neuralchecker/analyzing_constrained_LLM_through_PDFa_learning

A Proof of Proposition 2.1

Proof. Let u and v in Σ^* be arbitrary.

1. Assume that $u \equiv v$.

If $\mathbb{1}_{\mathcal{L}}(u) = 0$, then the lhs of (1) is undefined for any $w \in \Sigma^*$. Then $\mathbb{1}_{\mathcal{L}}(v) = 0$ since otherwise the rhs of (1) would be a number for any $w \in \Sigma^*$ (for instance it equals 1 for $w = \lambda$). By symmetry if $\mathbb{1}_{\mathcal{L}}(v) = 0$ then $\mathbb{1}_{\mathcal{L}}(u) = 0$. Therefore $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v)$.

Moreover, if $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v) = 0$ then $\mathbb{1}_{\mathcal{L}}(uw) = \mathbb{1}_{\mathcal{L}}(vw) = 0$ for all $w \in \Sigma^*$ and there is nothing more to check.

Suppose that $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v) = 1$ so that both sides of (1) are defined for any $w \in \Sigma^*$. Notice also that (1) implies $\mathbb{1}_{\mathcal{L}}(uw) = \mathbb{1}_{\mathcal{L}}(vw)$ for all $w \in \Sigma^*$. By definition of P we can rewrite (1) as follows:

$$\prod_{i=1}^{|w|} \mathcal{L}(u w_{<i})(w_i) = \prod_{i=1}^{|w|} \mathcal{L}(v w_{<i})(w_i) \quad (12)$$

for any $w \in \Sigma^*$ with length $|w| \geq 1$. In particular, varying $w = \sigma \in \Sigma$ in (12) and noticing that $\mathcal{L}(u)$ and $\mathcal{L}(v)$ are distributions over Σ_{\S} , we see that $\mathcal{L}(u) = \mathcal{L}(v)$.

We will now prove by induction on the length $|w|$ that $\mathcal{L}(uw) = \mathcal{L}(vw)$ whenever $\mathbb{1}_{\mathcal{L}}(uw) = \mathbb{1}_{\mathcal{L}}(vw) = 1$. We already proved the claim for $|w| = 0$, so suppose it holds true for length $\leq n$. Let w be such that $|w| = n + 1$ and let $\sigma \in \Sigma$ be such that $\mathbb{1}_{\mathcal{L}}(uw\sigma) = \mathbb{1}_{\mathcal{L}}(vw\sigma) = 1$. Since all terms involving the products in (12) are positive, and by induction hypothesis $\mathcal{L}(u w_{<i}) = \mathcal{L}(v w_{<i})$ for all $i = 1, \dots, n$, all these terms cancel out leaving the equality $\mathcal{L}(uw)(\sigma) = \mathcal{L}(vw)(\sigma)$. Since $\sigma \in \Sigma$ is arbitrary and $\mathcal{L}(uw)$ and $\mathcal{L}(vw)$ are probability distributions, we see again that they must be equal. This completes the proof.

2. Assume $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v)$ and $\forall w \in \Sigma^*. \mathbb{1}_{\mathcal{L}}(uw) = \mathbb{1}_{\mathcal{L}}(vw) = 1 \implies \mathcal{L}(uw) = \mathcal{L}(vw)$.

If $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v) = 0$, then the quotients in (1) are undefined and equality holds trivially for all $w \in \Sigma^*$.

Let us suppose then that $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v) = 1$. We first prove that $\mathbb{1}_{\mathcal{L}}(uw) = \mathbb{1}_{\mathcal{L}}(vw)$ for all $w \in \Sigma^*$. In fact, if on the contrary there exists $w \in \Sigma^*$ so that $\mathbb{1}_{\mathcal{L}}(uw) \neq \mathbb{1}_{\mathcal{L}}(vw)$, then there exists $w' \in \text{pref}(w)$ with $\mathbb{1}_{\mathcal{L}}(uw') = \mathbb{1}_{\mathcal{L}}(vw') = 1$ but $\mathcal{L}(uw') \neq \mathcal{L}(vw')$ because they have different support. This contradicts our assumption.

Let $w \in \Sigma^*$ be so that $\mathbb{1}_{\mathcal{L}}(uw) = \mathbb{1}_{\mathcal{L}}(vw) = 1$. Then for all prefix $w_{<i}$ we also have $\mathbb{1}_{\mathcal{L}}(uw_{<i}) = \mathbb{1}_{\mathcal{L}}(vw_{<i}) = 1$, and therefore $\mathcal{L}(uw_{<i}) = \mathcal{L}(vw_{<i})$. In particular, all the terms in (12) are equal and therefore (1) holds.

This completes the proof that $u \equiv v$.

□

B Proof of Proposition 2.2

Proof. Let $u \equiv_E v$. If $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v) = 0$, then $\mathbb{1}_{\mathcal{L}}(uw) = \mathbb{1}_{\mathcal{L}}(vw) = 0$ for all $w \in \Sigma^*$. Then $u\sigma \equiv_E v\sigma$ trivially.

Suppose now that $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v) = 1$ and let $\sigma \in \Sigma$. We have $\mathbb{1}_{\mathcal{L}}(u\sigma) = \mathbb{1}_{\mathcal{L}}(v\sigma)$ by Req. 3. Let $w \in \Sigma^*$ be arbitrary, since concatenation of strings is associative, if $\mathbb{1}_{\mathcal{L}}((u\sigma)w) = \mathbb{1}_{\mathcal{L}}((v\sigma)w) = 1$, then $\mathbb{1}_{\mathcal{L}}(u(\sigma w)) = \mathbb{1}_{\mathcal{L}}(v(\sigma w)) = 1$ and by assumption $\mathcal{L}(u(\sigma w)) =_E \mathcal{L}(v(\sigma w))$. Thus $\mathcal{L}((u\sigma)w) =_E \mathcal{L}((v\sigma)w)$. This proves that $u\sigma \equiv_E v\sigma$. □

C Proof of Proposition 2.3

Proof. Let $\alpha : \llbracket \Sigma^* \rrbracket_E \setminus \{0\} \rightarrow \Sigma^*$ be any function satisfying $\alpha(c) \in c$ for all $c \in \llbracket \Sigma^* \rrbracket_E \setminus \{0\}$. In other words, $\{\alpha(c) : c \in \llbracket \Sigma^* \rrbracket_E \setminus \{0\}\}$ is a set of representatives of the classes. Let $\beta : \Sigma^* \rightarrow \llbracket \Sigma^* \rrbracket_E^\bullet$ be the quotient map $\beta(u) = \llbracket u \rrbracket_E^\bullet$. Define $\phi = \beta \circ \alpha$.

Let $c, c' \in \llbracket \Sigma^* \rrbracket_E \setminus 0$ be such that $\phi(c) = \phi(c')$. Denote $u = \alpha(c)$ and $v = \alpha(c')$. By construction $\mathbb{1}_{\mathcal{L}}(u) = \mathbb{1}_{\mathcal{L}}(v) = 1$ and by Def. 5 we have $\mathcal{L}(uw) =_E \mathcal{L}(vw)$ for all $w \in \Sigma^*$. In particular $u \equiv_E v$, or equivalently $c = \llbracket u \rrbracket_E = \llbracket v \rrbracket_E = c'$. □

D Proof of Proposition 2.4

Proof. If $\mathbb{1}_{\mathcal{L}}(vw) = 1$ then $w' = w$. Otherwise, there exists $w'\sigma \in \text{pref}(w)$ such that $1 = \mathbb{1}_{\mathcal{L}}(vw') \neq \mathbb{1}_{\mathcal{L}}(vw'\sigma) = 0$. Hence, $\text{supp}(\mathcal{L}(uw')) \neq \text{supp}(\mathcal{L}(vw'))$ because $\mathbb{1}_{\mathcal{L}}(uw'\sigma) = 1$. Thus, by Req. 3, $\mathcal{L}(uw') \neq_E \mathcal{L}(vw')$. \square

E PDFA

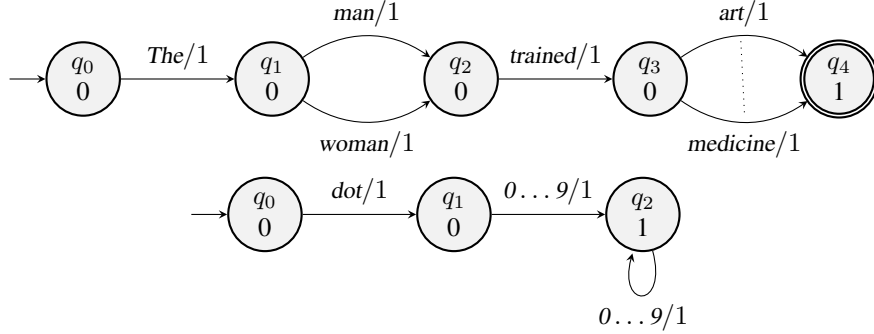


Figure 7: Guiding automata:(above) \mathcal{G}_1 for man-woman case study (below) \mathcal{G}_2 for digits case study

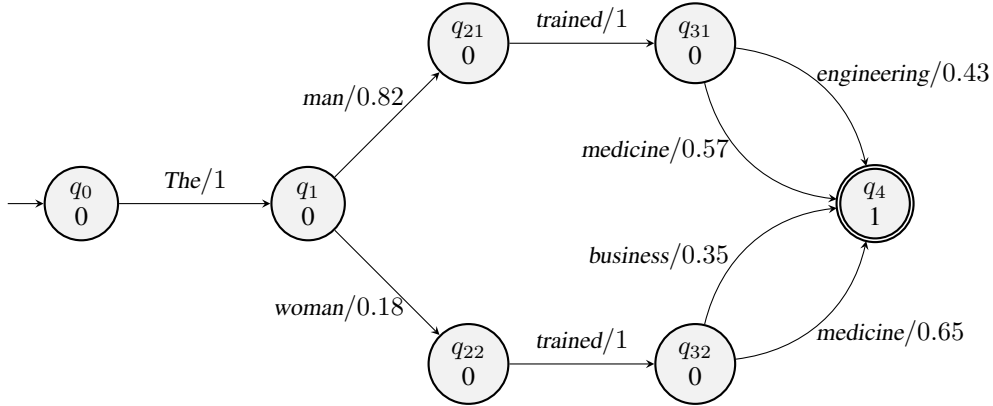


Figure 8: PDFA learnt for man-woman case study (with prefix space tokenizer)

F Existence of the probability measure P

Proposition F.1. *Let $\mathcal{L} : \Sigma^* \rightarrow \Delta(\Sigma_{\S})$ be a language model. There exists a unique probability measure P in $\Sigma^* \cup \Sigma^\omega$ such*

$$P(w) = \mathbf{P}\{x \in \Sigma^* \cup \Sigma^\omega : w \in \text{pref}(x)\} \text{ and } P_{\S}(w) = \mathbf{P}\{w\}$$

that for all $w \in \Sigma^*$.

Proof. We first extend the definition of \mathcal{L} in order to include the termination symbol. Let $\mathcal{L}_{\S} : \Sigma_{\S}^* \rightarrow \Delta(\Sigma_{\S})$ be defined as follows

$$\mathcal{L}_{\S}(w) = \begin{cases} \mathcal{L}(w) & \text{if } w \in \Sigma^* \\ \delta_{\S} & \text{if } w \in \Sigma_{\S}^* \setminus \Sigma^* \end{cases}$$

where $\delta_{\S}(\sigma) = 0$ for all $\sigma \in \Sigma$ and $\delta_{\S}(\S) = 1$. For each $k \geq 1$, we define the finite dimensional distribution $\mathbf{P}_k : \Sigma_{\S}^k \rightarrow [0, 1]$ as

$$\mathbf{P}_k[w] = \prod_{i=1}^k \mathcal{L}_{\S}(w_{<i})(w_i)$$

where we denote $w_{<i} = \sigma_1 \cdots \sigma_{i-1}$ if $w = \sigma_1 \cdots \sigma_k$, with the convention that $w_{<1} = \lambda$ the empty string. Let us show that $\{\mathbf{P}_k\}_{k \geq 1}$ is a consistent family of finite dimensional distributions:

$$\sum_{\sigma_{k+1}} \mathbf{P}_{k+1}(w\sigma_{k+1}) = \sum_{\sigma_{k+1}} \mathbf{P}_k[w] \mathcal{L}_{\S}(w)(\sigma_{k+1}) = \mathbf{P}_k[w] \sum_{\sigma_{k+1}} \mathcal{L}_{\S}(w)(\sigma_{k+1}) = \mathbf{P}_k[w]$$

By the Kolmogorov Extension Theorem (see [10] Thm.I.1.2) there exists a unique probability measure P in Σ_{\S}^ω such that $\mathbf{P}\{x \in \Sigma_{\S}^\omega : w \in \text{pref}(x)\} = \mathbf{P}_k[w]$ for all $k \geq 1$ and any $w \in \Sigma_{\S}^k$. Notice that in the usual measure theoretic terminology the event $\{x \in \Sigma_{\S}^\omega : w \in \text{pref}(x)\}$ is called a *cylinder*.

The event $A = \{x \in \Sigma_{\$}^{\omega} : \forall k \geq 1 \text{ if } x_k = \$ \text{ then } x_{k+1} = \$\}$ can be identified with $\Sigma^* \cup \Sigma^{\omega}$. Let us show that \mathbf{P} concentrates its measure in A , i.e. $\mathbf{P}[A] = 1$. The complement of A is

$$B = \bigcup_{k=1}^{\infty} B_k, \quad B_k = \{x \in \Sigma_{\$}^{\omega} : x_k = \$ \text{ and } x_{k+1} \neq \$\}$$

and B_k is the finite disjoint union of the cylinders of the form $C_{w,\sigma} = \{x \in \Sigma_{\$}^{\omega} : w\$ \sigma \in \text{Pref}(x)\}$ with $w \in \Sigma_{\$}^{k-1}$ and $\sigma \in \Sigma$. Therefore

$$\mathbf{P}[B_k] = \sum_{w,\sigma} \mathbf{P}[C_{w,\sigma}] = \sum_{w,\sigma} \mathbf{P}_{k+1}[C_{w,\sigma}] = \sum_{w,\sigma} \mathbf{P}_{k-1}[w] \mathcal{L}_{\$}(w)(\$) \delta_{\$}(\sigma) \overset{0}{=} 0$$

and the union bound shows that $\mathbf{P}[B] \leq \sum_{k=1}^{\infty} \mathbf{P}[B_k] = 0$.

Let us show the link between \mathbf{P} and P . Let us consider first a string $w \in \Sigma^*$ of length $k \geq 1$. Since the event $\{x \in \Sigma^* \cup \Sigma^{\omega} : w \in \text{pref}(x)\}$ equals the cylinder $C_k = \{x \in \Sigma_{\$}^{\omega} : w \in \text{pref}(x)\}$ intersected with A , and A has probability one, we have

$$\mathbf{P}\{x \in \Sigma^* \cup \Sigma^{\omega} : w \in \text{pref}(x)\} = \mathbf{P}_k[C_k] = \prod_{i=1}^k \mathcal{L}_{\$}(w_{<i})(w_i) = \prod_{i=1}^k \mathcal{L}(w_{<i})(w_i) = P(w)$$

In the case $w = \lambda$, the event $\{x \in \Sigma^* \cup \Sigma^{\omega} : w \in \text{pref}(x)\}$ equals A and its probability is therefore 1 as it is the case for $P(\lambda)$.

Finally, let us compute the probability of occurrence of a given finite string $w \in \Sigma^*$. This string corresponds to the infinite sequence $w\$\$\$\dots$ in $\Sigma_{\$}^{\omega}$, which in turn equals the decreasing intersection of the cylinders $C_{w,n} = \{x \in \Sigma_{\$}^{\omega} : w(\$)^n \in \text{pref}(x)\}$. Therefore

$$\mathbf{P}\{w\} = \mathbf{P}\left[\bigcap_{n \geq 1} C_{w,n}\right] = \lim_{n \rightarrow +\infty} \left[\prod_{i=1}^{|w|} \mathcal{L}(w_{<i})(w_i) \right] \mathcal{L}(w)(\$) \left[\prod_{j=0}^{n-1} \delta_{\$}(\$) \right] \overset{1}{=} P_{\$}(w)$$

This concludes the proof. □