

MAXIMUM LIKELIHOOD ESTIMATION OF FACTORED REGULAR DETERMINISTIC STOCHASTIC LANGUAGES

Chihiro Shibata and Jeffrey Heinz



University of Toronto
July 19, 2019

We thank JSPS KAKENHI #JP18K11449 (CS) and
NIH #R01HD87133-01 (JH)

THE PROBLEM IN GENERAL

Stochastic languages are probability distributions over strings.

THE PROBLEM IN GENERAL

Stochastic languages are probability distributions over strings.

$$f : \Sigma^* \rightarrow [0, 1] \text{ and } \sum_w f(w) = 1$$

THE PROBLEM IN GENERAL

Stochastic languages are probability distributions over strings.

$$f : \Sigma^* \rightarrow [0, 1] \text{ and } \sum_w f(w) = 1$$

A class C of stochastic languages is often defined parametrically: an assignment of values to parameters Θ uniquely determines some stochastic language $f_\Theta \in C$.

THE PROBLEM IN GENERAL

Stochastic languages are probability distributions over strings.

$$f : \Sigma^* \rightarrow [0, 1] \text{ and } \sum_w f(w) = 1$$

A class C of stochastic languages is often defined parametrically: an assignment of values to parameters Θ uniquely determines some stochastic language $f_\Theta \in C$.

An important learning criterion

For *any* data sequence D drawn i.i.d. from *any* stochastic language, a Maximum-Likelihood Estimator finds parameter values $\hat{\Theta}$ which maximize $P(D)$ with respect to C .

THE PROBLEM IN GENERAL

Stochastic languages are probability distributions over strings.

$$f : \Sigma^* \rightarrow [0, 1] \text{ and } \sum_w f(w) = 1$$

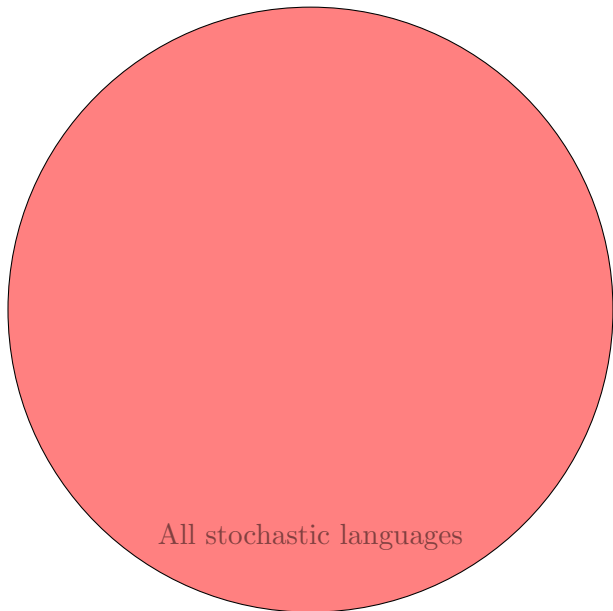
A class C of stochastic languages is often defined parametrically: an assignment of values to parameters Θ uniquely determines some stochastic language $f_\Theta \in C$.

An important learning criterion

For *any* data sequence D drawn i.i.d. from *any* stochastic language, a Maximum-Likelihood Estimator finds parameter values $\hat{\Theta}$ which maximize $P(D)$ with respect to C .

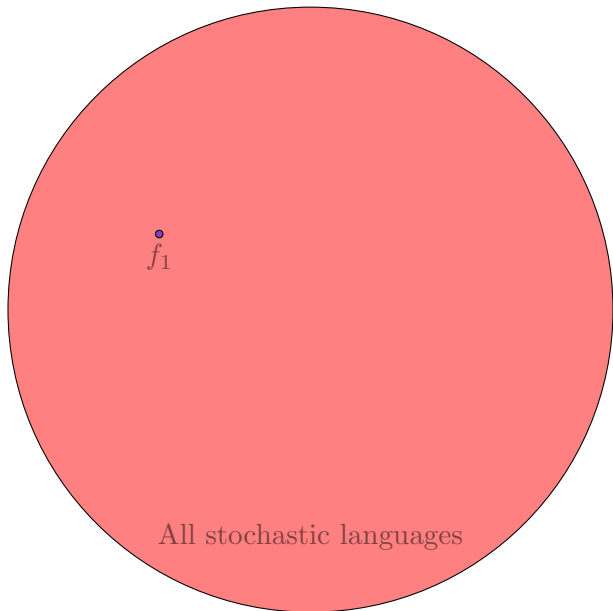
For a class of stochastic languages C , is there an algorithm which reliably returns a Maximum-Likelihood Estimate (MLE) of an observed data sample D ?

IN PICTURES

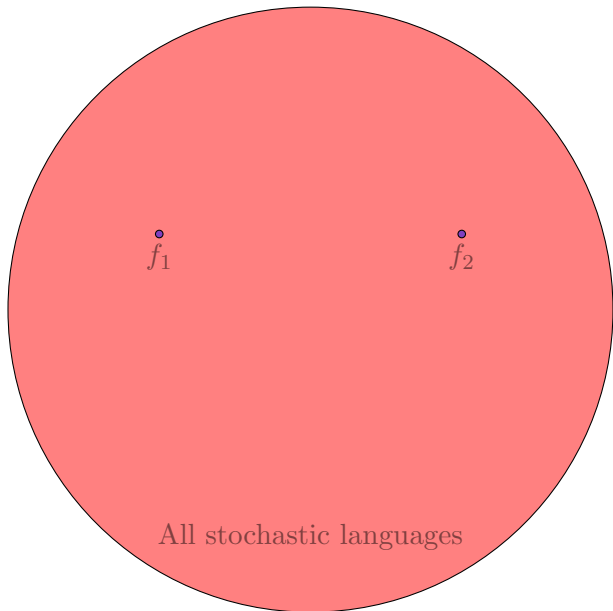


All stochastic languages

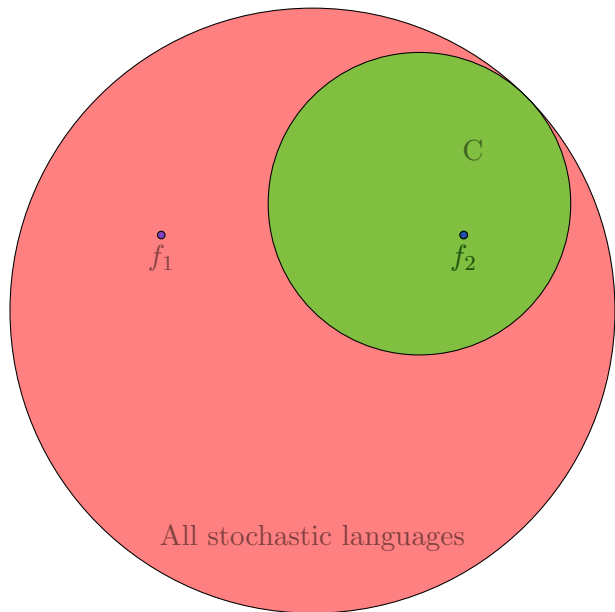
IN PICTURES



IN PICTURES



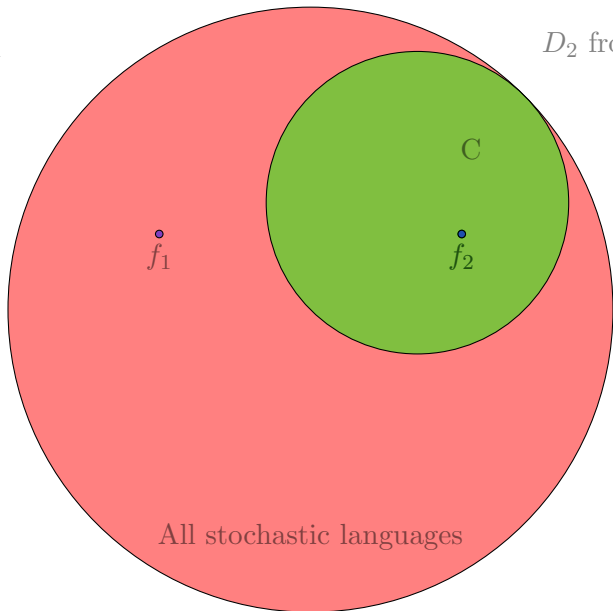
IN PICTURES



IN PICTURES

D_1 from f_1

D_2 from f_2



f_1

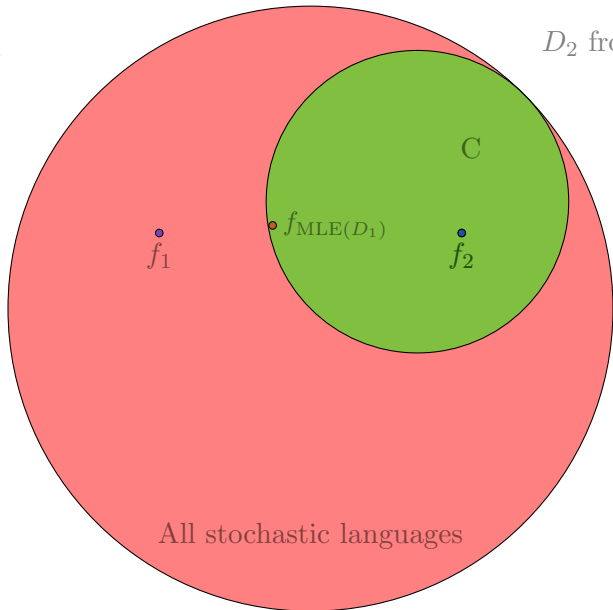
f_2

All stochastic languages

IN PICTURES

D_1 from f_1

D_2 from f_2



f_1

$f_{MLE(D_1)}$

f_2

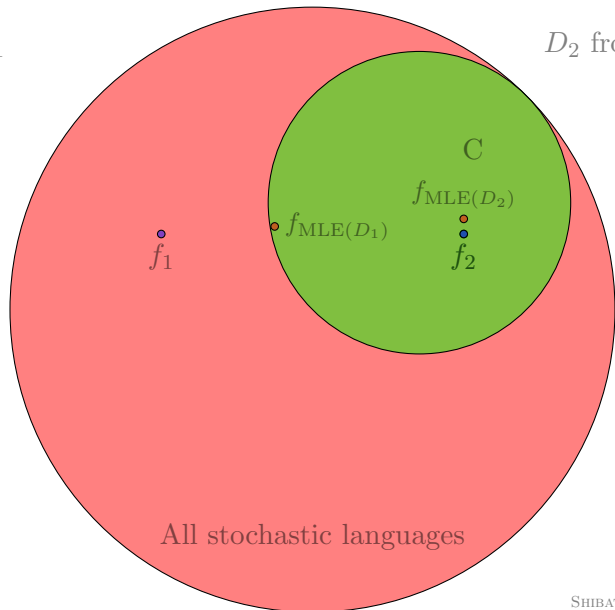
C

All stochastic languages

IN PICTURES

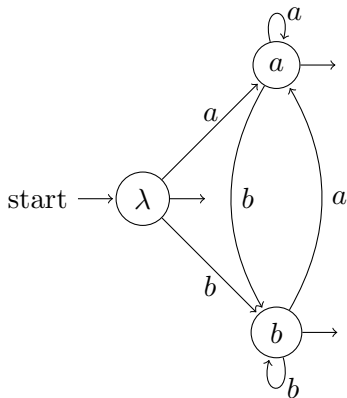
D_1 from f_1

D_2 from f_2



CLASSES DEFINED BY SINGLE DFAS

Example: Bigram model



Parameters

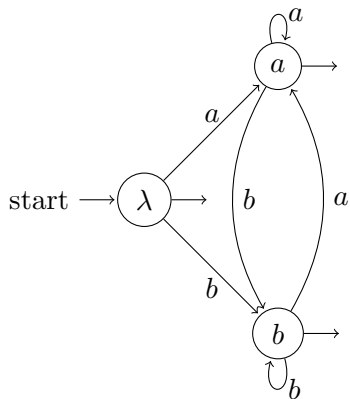
 $\theta_{\times a}$ $\theta_{\times b}$ $\theta_{\times \times}$

 θ_{aa} θ_{ab} $\theta_{a \times}$

 θ_{ba} θ_{bb} $\theta_{b \times}$

CLASSES DEFINED BY SINGLE DFAS

Example: Bigram model



Parameters

 $\theta_{\times a}$ $\theta_{\times b}$ $\theta_{\times \times}$

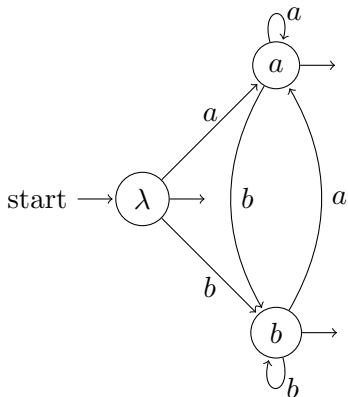
 θ_{aa} θ_{ab} $\theta_{a \times}$

 θ_{ba} θ_{bb} $\theta_{b \times}$

$$D = \langle ab, aabb \rangle$$

CLASSES DEFINED BY SINGLE DFAS

Example: Bigram model



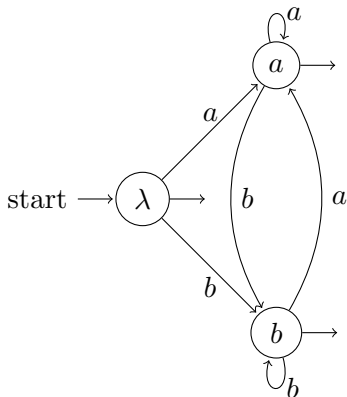
Parameters	
$\theta_{\times a}$	1
$\theta_{\times b}$	
$\theta_{\times \times}$	
θ_{aa}	
θ_{ab}	1
$\theta_{a \times}$	
θ_{ba}	
θ_{bb}	
$\theta_{b \times}$	1

$$D = \langle ab, aabb \rangle$$

MLE is obtained by passing D through DFA and normalizing.
(Vidal et al. 2005)

CLASSES DEFINED BY SINGLE DFAS

Example: Bigram model



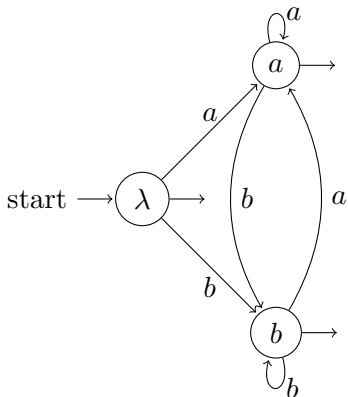
Parameters	
$\theta_{\times a}$	2
$\theta_{\times b}$	
$\theta_{\times \times}$	
θ_{aa}	1
θ_{ab}	2
$\theta_{a \times}$	
θ_{ba}	
θ_{bb}	1
$\theta_{b \times}$	2

$$D = \langle ab, aabb \rangle$$

MLE is obtained by passing D through DFA and normalizing.
(Vidal et al. 2005)

CLASSES DEFINED BY SINGLE DFAS

Example: Bigram model



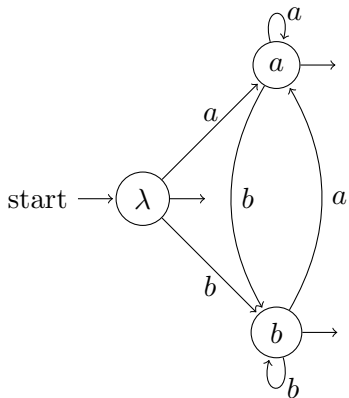
Parameters	
$\theta_{\times a}$	1
$\theta_{\times b}$	0
$\theta_{\times \times}$	0
θ_{aa}	1/3
θ_{ab}	2/3
$\theta_{a \times}$	0
θ_{ba}	0
θ_{bb}	1/3
$\theta_{b \times}$	2/3

$$D = \langle ab, aabb \rangle$$

MLE is obtained by passing D through DFA and normalizing.
(Vidal et al. 2005)

CLASSES DEFINED BY SINGLE DFAS

Example: Strictly 2-Local Languages



Parameters

$\theta_{\times a}$

$\theta_{\times b}$

$\theta_{\times \times}$

θ_{aa}

θ_{ab}

$\theta_{a \times}$

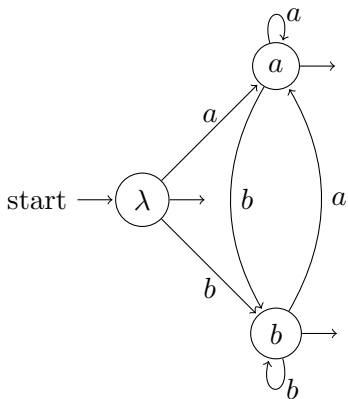
θ_{ba}

θ_{bb}

$\theta_{b \times}$

CLASSES DEFINED BY SINGLE DFAS

Example: Strictly 2-Local Languages



Parameters

$\theta_{\times a}$

$\theta_{\times b}$

$\theta_{\times \times}$

θ_{aa}

θ_{ab}

$\theta_{a \times}$

θ_{ba}

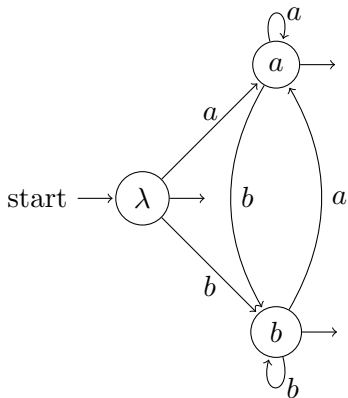
θ_{bb}

$\theta_{b \times}$

$$D = \langle ab, aabb \rangle$$

CLASSES DEFINED BY SINGLE DFAS

Example: Strictly 2-Local Languages



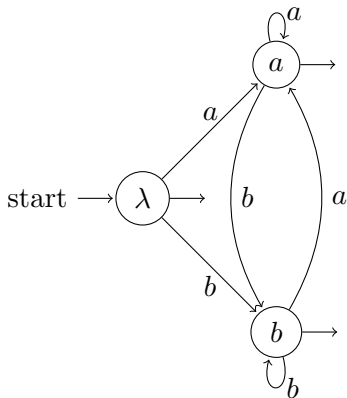
Parameters	
$\theta_{\times a}$	1
$\theta_{\times b}$	
$\theta_{\times \times}$	
θ_{aa}	
θ_{ab}	1
$\theta_{a \times}$	
θ_{ba}	
θ_{bb}	
$\theta_{b \times}$	1

$$D = \langle ab, aabb \rangle$$

Smallest language consistent with D in C is obtained by passing D through DFA and ‘activating’ parsed transitions. (Heinz and Rogers 2013)

CLASSES DEFINED BY SINGLE DFAS

Example: Strictly 2-Local Languages



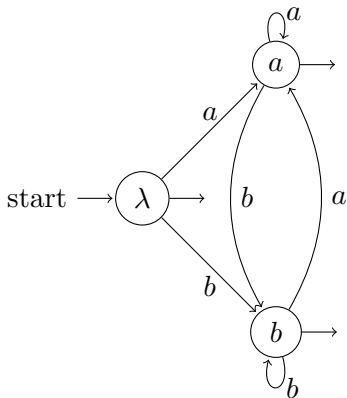
Parameters	
$\theta_{\times a}$	1
$\theta_{\times b}$	
$\theta_{\times \times}$	
θ_{aa}	1
θ_{ab}	1
$\theta_{a \times}$	
θ_{ba}	
θ_{bb}	1
$\theta_{b \times}$	1

$$D = \langle ab, aabb \rangle$$

Smallest language consistent with D in C is obtained by passing D through DFA and ‘activating’ parsed transitions. (Heinz and Rogers 2013)

CLASSES DEFINED BY SINGLE DFAS

Example: Strictly 2-Local Languages



Parameters	
$\theta_{\times a}$	1
$\theta_{\times b}$	0
$\theta_{\times \times}$	0
θ_{aa}	1
θ_{ab}	1
$\theta_{a \times}$	0
θ_{ba}	0
θ_{bb}	1
$\theta_{b \times}$	1

$$D = \langle ab, aabb \rangle$$

Smallest language consistent with D in C is obtained by passing D through DFA and ‘activating’ parsed transitions. (Heinz and Rogers 2013)

OVERVIEW OF RELATED RESULTS

		Class C defined with	
		single DFA	finitely many DFA
type of	$f : \Sigma^* \rightarrow \{0, 1\}$		
language	$f : \Sigma^* \rightarrow [0, 1]$		

- 1 For Boolean languages, the learning algorithms return the smallest language in C which includes D .
- 2 For Stochastic languages, the MLE returns the language in C which maximizes likelihood of D .

(Vidal et al. 2005, Heinz and Rogers 2013)

OVERVIEW OF RELATED RESULTS

		Class C defined with	
		single DFA	finitely many DFA
type of	$f : \Sigma^* \rightarrow \{0, 1\}$	✓	
language	$f : \Sigma^* \rightarrow [0, 1]$		

- 1 For Boolean languages, the learning algorithms return the smallest language in C which includes D .
- 2 For Stochastic languages, the MLE returns the language in C which maximizes likelihood of D .

(Vidal et al. 2005, Heinz and Rogers 2013)

OVERVIEW OF RELATED RESULTS

		Class C defined with	
		single DFA	finitely many DFA
type of	$f : \Sigma^* \rightarrow \{0, 1\}$	✓	
language	$f : \Sigma^* \rightarrow [0, 1]$	✓	

- 1 For Boolean languages, the learning algorithms return the smallest language in C which includes D .
- 2 For Stochastic languages, the MLE returns the language in C which maximizes likelihood of D .

(Vidal et al. 2005, Heinz and Rogers 2013)

OVERVIEW OF RELATED RESULTS

		Class C defined with	
		single DFA	finitely many DFA
type of	$f : \Sigma^* \rightarrow \{0, 1\}$	✓	✓
language	$f : \Sigma^* \rightarrow [0, 1]$	✓	

- 1 For Boolean languages, the learning algorithms return the smallest language in C which includes D .
- 2 For Stochastic languages, the MLE returns the language in C which maximizes likelihood of D .

(Vidal et al. 2005, Heinz and Rogers 2013)

OVERVIEW OF RELATED RESULTS

		Class C defined with	
		single DFA	finitely many DFA
type of	$f : \Sigma^* \rightarrow \{0, 1\}$	✓	✓
language	$f : \Sigma^* \rightarrow [0, 1]$	✓	this paper

- 1 For Boolean languages, the learning algorithms return the smallest language in C which includes D .
- 2 For Stochastic languages, the MLE returns the language in C which maximizes likelihood of D .

(Vidal et al. 2005, Heinz and Rogers 2013)

OVERVIEW OF RELATED RESULTS (PART 2)

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).

OVERVIEW OF RELATED RESULTS (PART 2)

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).
- 2 It is NP-hard to find the minimal DFA consistent with a finite sample of positive and negative examples (Gold 1978).

OVERVIEW OF RELATED RESULTS (PART 2)

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).
- 2 It is NP-hard to find the minimal DFA consistent with a finite sample of positive and negative examples (Gold 1978).
- 3 Each DFA admits a characteristic sample D of positive and negative examples such that RPNI identifies the DFA from any superset of D in cubic time (Oncina and Garcia 1992, DuPont 1996).

OVERVIEW OF RELATED RESULTS (PART 2)

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).
- 2 It is NP-hard to find the minimal DFA consistent with a finite sample of positive and negative examples (Gold 1978).
- 3 Each DFA admits a characteristic sample D of positive and negative examples such that RPNI identifies the DFA from any superset of D in cubic time (Oncina and Garcia 1992, DuPont 1996).
- 4 ALEGRIA/RLIPS (based on RPNI) (Carrasco and Oncina 1994, 1999) learns the class of PDFAs in polynomial time with probability one (de la Higuera and Thollard 2001).

OVERVIEW OF RELATED RESULTS (PART 2)

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).
- 2 It is NP-hard to find the minimal DFA consistent with a finite sample of positive and negative examples (Gold 1978).
- 3 Each DFA admits a characteristic sample D of positive and negative examples such that RPNI identifies the DFA from any superset of D in cubic time (Oncina and Garcia 1992, DuPont 1996).
- 4 ALEGRIA/RLIPS (based on RPNI) (Carrasco and Oncina 1994, 1999) learns the class of PDFAs in polynomial time with probability one (de la Higuera and Thollard 2001).
- 5 Clark and Thollard (2004) present an algorithm which learns the class of PDFAs in a modified PAC setting. (See also Parekh and Hanover 2001.)

OVERVIEW OF RELATED RESULTS (PART 2)

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).
- 2 It is NP-hard to find the minimal DFA consistent with a finite sample of positive and negative examples (Gold 1978).
- 3 Each DFA admits a characteristic sample D of positive and negative examples such that RPNI identifies the DFA from any superset of D in cubic time (Oncina and Garcia 1992, DuPont 1996).
- 4 ALEGRIA/RLIPS (based on RPNI) (Carrasco and Oncina 1994, 1999) learns the class of PDFAs in polynomial time with probability one (de la Higuera and Thollard 2001).
- 5 Clark and Thollard (2004) present an algorithm which learns the class of PDFAs in a modified PAC setting. (See also Parekh and Hanover 2001.)
- 6 Maximization-Expectation techniques are used to learn the class of PNFAs, but there is no guarantee to find a global optimum (Rabiner 1989).

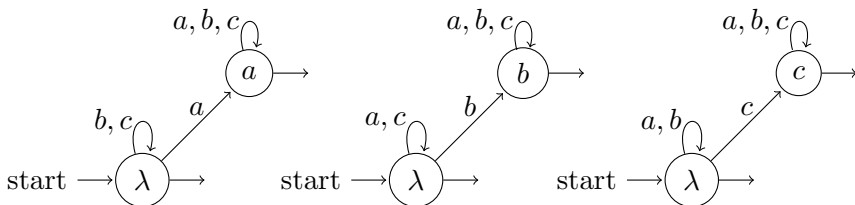
DEFINING C WITH FINITELY MANY DFA

DEFINING C WITH FINITELY MANY DFA

How do you define a class C with finitely many DFA?

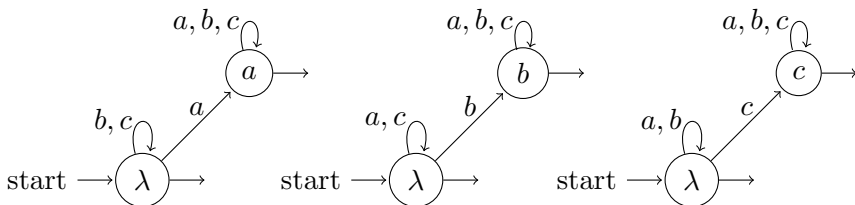
DEFINING C WITH FINITELY MANY DFA

How do you define a class C with finitely many DFA?



DEFINING C WITH FINITELY MANY DFA

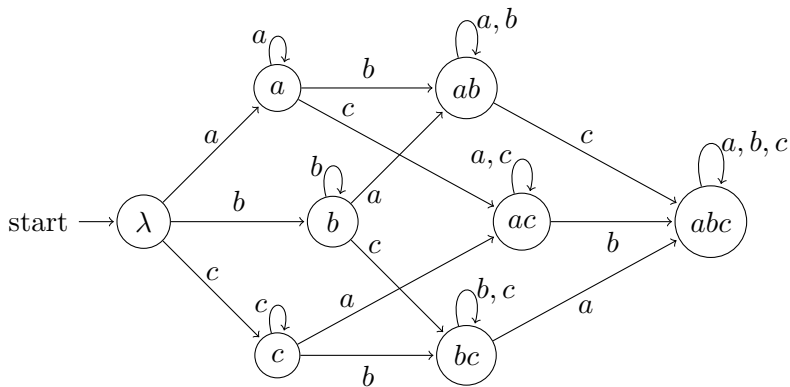
How do you define a class C with finitely many DFA?



Product Operations

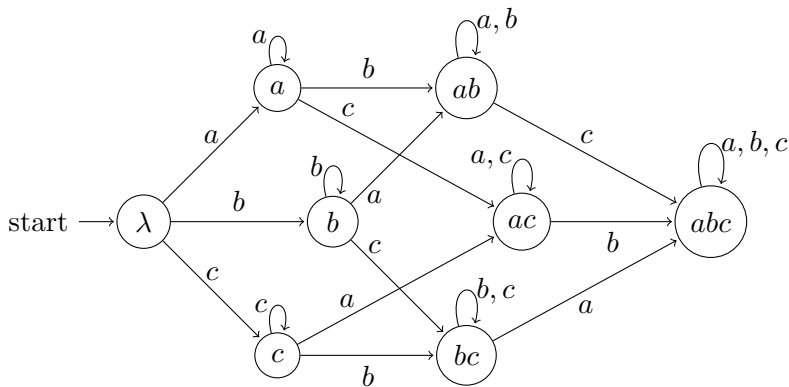
- 1 For Boolean languages, use **acceptor product** (yields intersection)
- 2 For Stochastic languages, use **co-emission product** (yields joint distribution)

THE PRODUCT OF THOSE THREE ACCEPTORS



(exit/accepting arrow at each state is not shown)

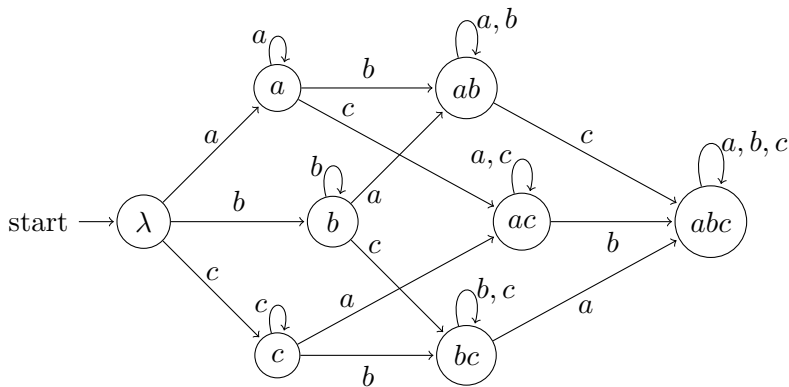
THE PRODUCT OF THOSE THREE ACCEPTORS



(exit/accepting arrow at each state is not shown)

- If C is defined by this DFA, then $C = \text{Piecewise 2-Testable}$.

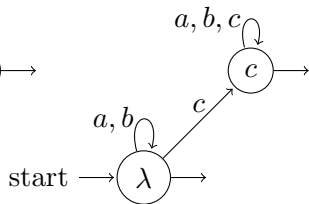
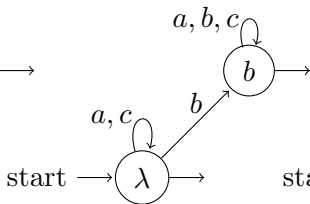
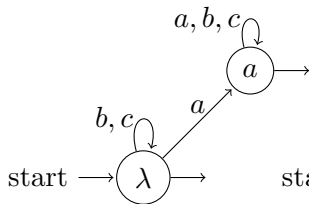
THE PRODUCT OF THOSE THREE ACCEPTORS



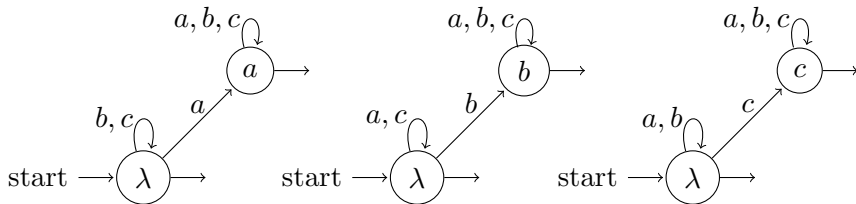
(exit/accepting arrow at each state is not shown)

- If C is defined by this DFA, then $C = \text{Piecewise 2-Testable}$.
- If C is defined by the 3 atomic DFAs, then $C = \text{Strictly 2-Piecewise}$.

CAUSE . . .

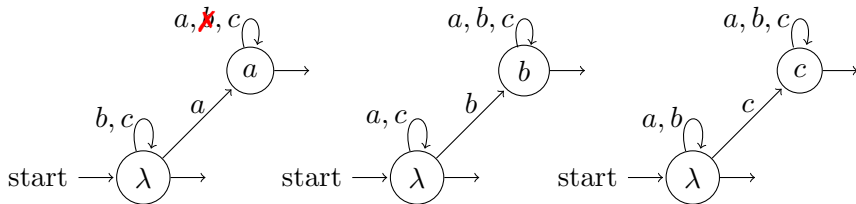


CAUSE . . .



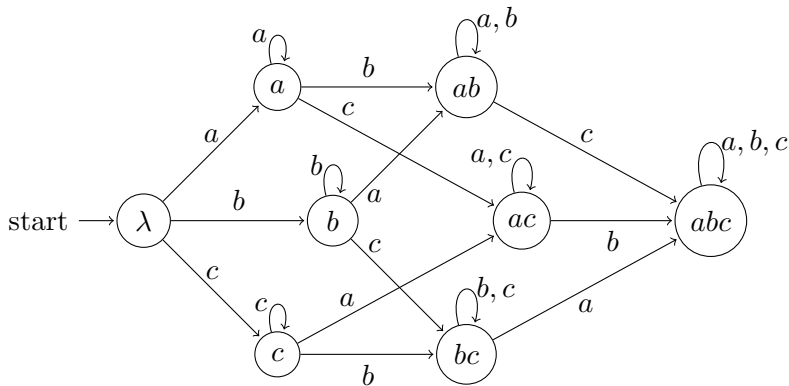
The parameters of the model are set at the level of the individual DFA.

CAUSE . . .



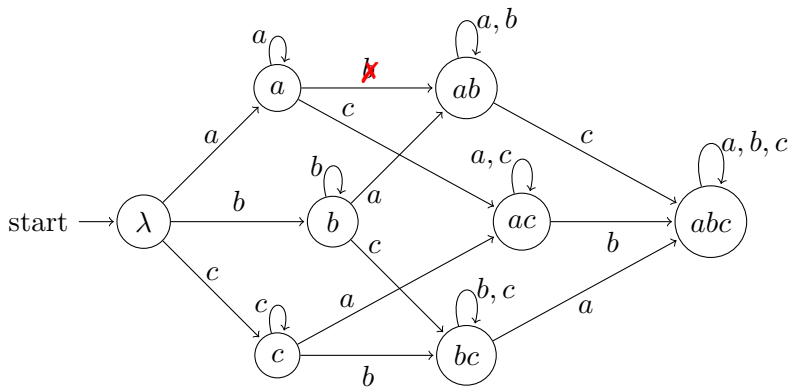
The parameters of the model are set at the level of the individual DFA.

... AND EFFECT



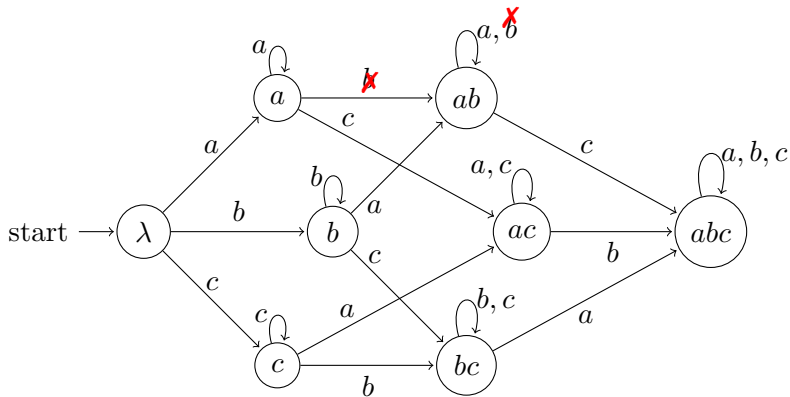
(exit/accepting arrow at each state is not shown)

... AND EFFECT



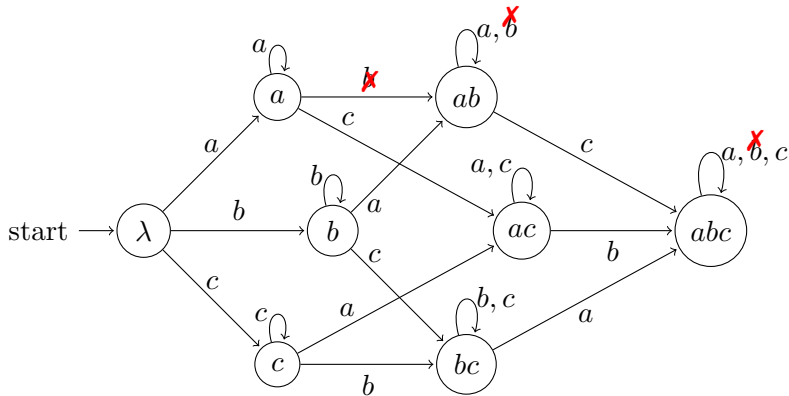
(exit/accepting arrow at each state is not shown)

... AND EFFECT



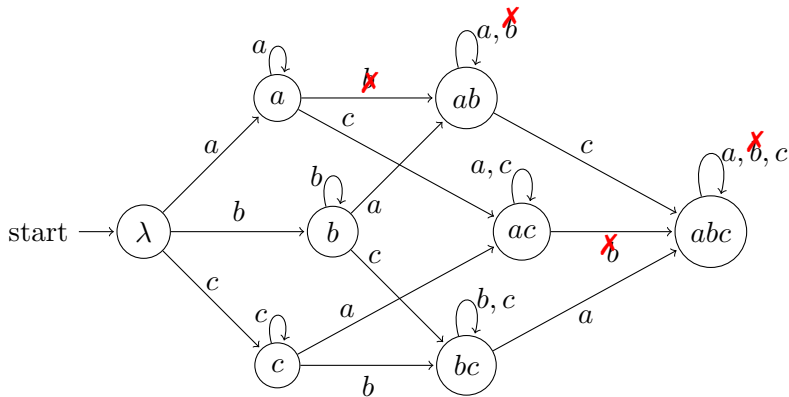
(exit/accepting arrow at each state is not shown)

... AND EFFECT



(exit/accepting arrow at each state is not shown)

... AND EFFECT



(exit/accepting arrow at each state is not shown)

COMPARING THE REPRESENTATIONS

The Product DFA

The Atomic DFAs

COMPARING THE REPRESENTATIONS

The Product DFA

- 1 In the worst case, it has $\prod_i |Q_i|$ states and $(|\Sigma| + 1) \prod_i |Q_i|$ parameters.

The Atomic DFAs

COMPARING THE REPRESENTATIONS

The Product DFA

- 1 In the worst case, it has $\prod_i |Q_i|$ states and $(|\Sigma| + 1) \prod_i |Q_i|$ parameters.

The Atomic DFAs

- 1 The atomic DFAs have a total of $\sum_i |Q_i|$ states and $(|\Sigma| + 1) \sum_i |Q_i|$ parameters.

COMPARING THE REPRESENTATIONS

The Product DFA

- 1 In the worst case, it has $\prod_i |Q_i|$ states and $(|\Sigma| + 1) \prod_i |Q_i|$ parameters.
- 2 Transitions/parameters are independent of others.

The Atomic DFAs

- 1 The atomic DFAs have a total of $\sum_i |Q_i|$ states and $(|\Sigma| + 1) \sum_i |Q_i|$ parameters.

COMPARING THE REPRESENTATIONS

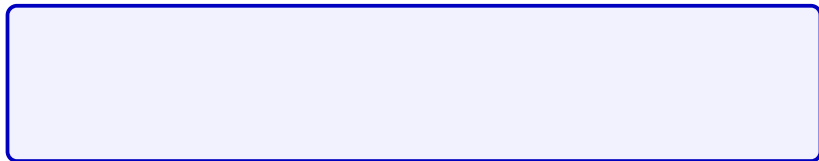
The Product DFA

- 1 In the worst case, it has $\prod_i |Q_i|$ states and $(|\Sigma| + 1) \prod_i |Q_i|$ parameters.
- 2 Transitions/parameters are independent of others.

The Atomic DFAs

- 1 The atomic DFAs have a total of $\sum_i |Q_i|$ states and $(|\Sigma| + 1) \sum_i |Q_i|$ parameters.
- 2 The transitions in the product are NOT independent.

PLUSES AND MINUSES



PLUSES AND MINUSES

+ Fewer parameters means more accurate estimation of model parameters with less data.

PLUSES AND MINUSES

- + Fewer parameters means more accurate estimation of model parameters with less data.
- Fewer parameters means the model is less expressive.

PLUSES AND MINUSES

- + Fewer parameters means more accurate estimation of model parameters with less data.
- Fewer parameters means the model is less expressive.

- Heinz and Rogers (2013, MoL) extend the method of ‘activating’ data-parsed transitions to learn classes of Boolean languages defined with single DFA to classes of Boolean languages defined with finitely many DFA.

PLUSES AND MINUSES

- + Fewer parameters means more accurate estimation of model parameters with less data.
- Fewer parameters means the model is less expressive.

- Heinz and Rogers (2013, MoL) extend the method of ‘activating’ data-parsed transitions to learn classes of Boolean languages defined with single DFA to classes of Boolean languages defined with finitely many DFA.
- They show it always returns the smallest Boolean language in the class consistent with the data, and thus identifies the class in the limit from positive data.

THE CO-EMISSION PRODUCT

THE CO-EMISSION PRODUCT

- The co-emission product defines how PDFA-definable stochastic languages can be multiplied together to yield a well-defined stochastic language.

THE CO-EMISSION PRODUCT

- The co-emission product defines how PDFFA-definable stochastic languages can be multiplied together to yield a well-defined stochastic language.
- Heinz and Rogers 2010 defined stochastic Strictly k -Piecewise languages using a *variant* of the co-emission product.

THE CO-EMISSION PRODUCT

- The co-emission product defines how PDFFA-definable stochastic languages can be multiplied together to yield a well-defined stochastic language.
- Heinz and Rogers 2010 defined stochastic Strictly k -Piecewise languages using a *variant* of the co-emission product.
- They claimed they could find the MLE, but nobody seemed convinced.

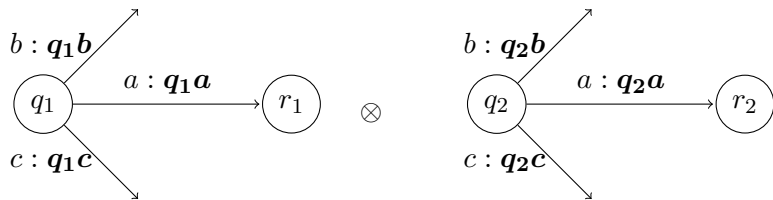
THE CO-EMISSION PRODUCT

- The co-emission product defines how PDFA-definable stochastic languages can be multiplied together to yield a well-defined stochastic language.
- Heinz and Rogers 2010 defined stochastic Strictly k -Piecewise languages using a *variant* of the co-emission product.
- They claimed they could find the MLE, but nobody seemed convinced.

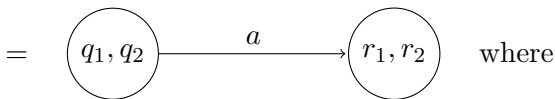
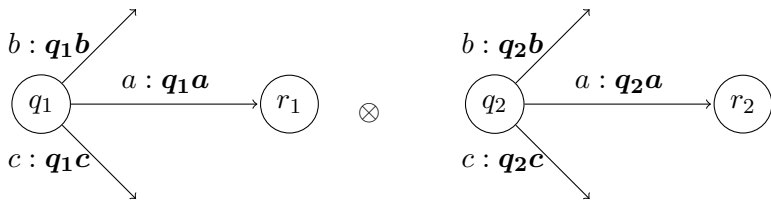
$Pr(x \mid P_{\leq 1}(y))$		x			
		s	\widehat{ts}	\int	\widehat{tj}
y	s	0.0335	0.0051	0.0011	0.0002
	\widehat{ts}	0.0218	0.0113	0.0009	0.
	\int	0.0009	0.	0.0671	0.0353
	\widehat{tj}	0.0006	0.	0.0455	0.0313

TABLE: Results of SP_2 estimation on the Samala corpus. Only sibilants are shown. (Heinz and Rogers 2010, p. 894)

THE CO-EMISSION PRODUCT (DEFINITION)

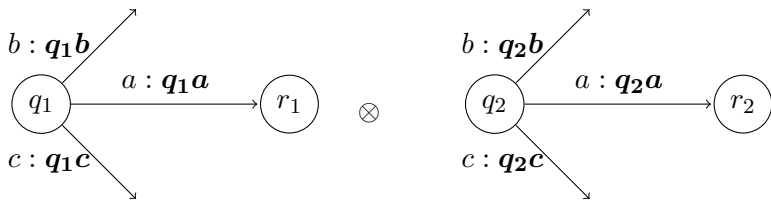


THE CO-EMISSION PRODUCT (DEFINITION)



$$P(a \mid q_1, q_2) \stackrel{\text{def}}{=} \frac{\prod_i q_i a}{\sum_{\sigma} \prod_i q_i \sigma}$$

THE CO-EMISSION PRODUCT (DEFINITION)



$$= \text{node}(q_1, q_2) \xrightarrow{a} \text{node}(r_1, r_2) \quad \text{where}$$

$$P(a \mid q_1, q_2) \stackrel{\text{def}}{=} \frac{\prod_i q_i a}{\sum_{\sigma} \prod_i q_i \sigma}$$

For fixed σ , the co-emission product treats the parameters $q_i \sigma$ as *independent*.

CONTRIBUTIONS

- 1 We extend Heinz and Rogers 2010 analysis to classes defined with

CONTRIBUTIONS

- 1 We extend Heinz and Rogers 2010 analysis to classes defined with
 - 1 the standard co-emission product (not the variant introduced by Heinz and Rogers)

CONTRIBUTIONS

- 1 We extend Heinz and Rogers 2010 analysis to classes defined with
 - 1 the standard co-emission product (not the variant introduced by Heinz and Rogers)
 - 2 of arbitrary sets of finitely many PDFAs (not just the ones which define stochastic SP_k languages)

CONTRIBUTIONS

- ① We extend Heinz and Rogers 2010 analysis to classes defined with
 - ① the standard co-emission product (not the variant introduced by Heinz and Rogers)
 - ② of arbitrary sets of finitely many PDFAs (not just the ones which define stochastic SP_k languages)
- ② Essentially, we prove that parameters which maximize the probability of the data with respect to such models are found by running the corpus through each of the individual factor PDFAs and calculating the relative frequencies.

SOME DETAILS OF THE ANALYSIS

- 1 Probability of Words
- 2 Relative Frequency of Emissions
- 3 Empirical Mean of co-emission probabilities
- 4 Main Theorems

PROBABILITY OF WORDS

- Consider a class C defined with the co-emission product of K machines $M_1 \dots M_K$.

PROBABILITY OF WORDS

- Consider a class C defined with the co-emission product of K machines $M_1 \dots M_K$.
- Suppose that $w = \sigma_1 \dots \sigma_N$

PROBABILITY OF WORDS

- Consider a class C defined with the co-emission product of K machines $M_1 \dots M_K$.
- Suppose that $w = \sigma_1 \dots \sigma_N$
- Let $q(j, i)$ denote a state in Q_j that is reached after M_j reads the prefix $\sigma_1 \dots \sigma_{i-1}$.

PROBABILITY OF WORDS

- Consider a class C defined with the co-emission product of K machines $M_1 \dots M_K$.
- Suppose that $w = \sigma_1 \dots \sigma_N$
- Let $q(j, i)$ denote a state in Q_j that is reached after M_j reads the prefix $\sigma_1 \dots \sigma_{i-1}$.
- If $i = 1$ then $q(j, i)$ represents the initial state of M_j .

PROBABILITY OF WORDS

- Consider a class C defined with the co-emission product of K machines $M_1 \dots M_K$.
- Suppose that $w = \sigma_1 \dots \sigma_N$
- Let $q(j, i)$ denote a state in Q_j that is reached after M_j reads the prefix $\sigma_1 \dots \sigma_{i-1}$.
- If $i = 1$ then $q(j, i)$ represents the initial state of M_j .
- Let $T_j(q, \sigma)$ denote a parameter (transitional probability) in PDFFA M_j .

PROBABILITY OF WORDS

- Consider a class C defined with the co-emission product of K machines $M_1 \dots M_K$.
- Suppose that $w = \sigma_1 \dots \sigma_N$
- Let $q(j, i)$ denote a state in Q_j that is reached after M_j reads the prefix $\sigma_1 \dots \sigma_{i-1}$.
- If $i = 1$ then $q(j, i)$ represents the initial state of M_j .
- Let $T_j(q, \sigma)$ denote a parameter (transitional probability) in PDFA M_j .
- Then the probability that σ is emitted after the product machine $\bigotimes_{1 \leq j \leq K} \mathcal{M}_j$ reads the prefix $\sigma_1 \dots \sigma_{i-1}$ is the following:

$$\text{Coemit}(\sigma, i) = \frac{\prod_{j=1}^K T_j(q(j, i), \sigma)}{\sum_{\sigma' \in \Sigma} \prod_{j=1}^K T_j(q(j, i), \sigma')}. \quad (1)$$

PROBABILITY OF WORDS

- Consider a class C defined with the co-emission product of K machines $M_1 \dots M_K$.
- Suppose that $w = \sigma_1 \dots \sigma_N$
- Let $q(j, i)$ denote a state in Q_j that is reached after M_j reads the prefix $\sigma_1 \dots \sigma_{i-1}$.
- If $i = 1$ then $q(j, i)$ represents the initial state of M_j .
- Let $T_j(q, \sigma)$ denote a parameter (transitional probability) in PDFAs M_j .
- Then the probability that σ is emitted after the product machine $\bigotimes_{1 \leq j \leq K} \mathcal{M}_j$ reads the prefix $\sigma_1 \dots \sigma_{i-1}$ is the following:

$$\text{Coemit}(\sigma, i) = \frac{\prod_{j=1}^K T_j(q(j, i), \sigma)}{\sum_{\sigma' \in \Sigma} \prod_{j=1}^K T_j(q(j, i), \sigma')}. \quad (1)$$

- We assume that there is an end marker $\times \in \Sigma$ which uniquely occurs at the end of words.

PROBABILITY OF WORDS

- Consider a class C defined with the co-emission product of K machines $M_1 \dots M_K$.
- Suppose that $w = \sigma_1 \dots \sigma_N$
- Let $q(j, i)$ denote a state in Q_j that is reached after M_j reads the prefix $\sigma_1 \dots \sigma_{i-1}$.
- If $i = 1$ then $q(j, i)$ represents the initial state of M_j .
- Let $T_j(q, \sigma)$ denote a parameter (transitional probability) in PDFA M_j .
- Then the probability that σ is emitted after the product machine $\bigotimes_{1 \leq j \leq K} \mathcal{M}_j$ reads the prefix $\sigma_1 \dots \sigma_{i-1}$ is the following:

$$\text{Coemit}(\sigma, i) = \frac{\prod_{j=1}^K T_j(q(j, i), \sigma)}{\sum_{\sigma' \in \Sigma} \prod_{j=1}^K T_j(q(j, i), \sigma')}. \quad (1)$$

- We assume that there is a end marker $\times \in \Sigma$ which uniquely occurs at the end of words.

$$P(w \times) = \prod_{i=1}^{N+1} \text{Coemit}(\sigma_i, i) \quad (2)$$

RELATIVE FREQUENCY OF EMISSION

RELATIVE FREQUENCY OF EMISSION

- Let $m_w(M_j, q, \sigma) \in \mathbb{Z}^+$ denote how many times σ is emitted at the state q while the machine M_j emits w .

RELATIVE FREQUENCY OF EMISSION

- Let $m_w(M_j, q, \sigma) \in \mathbb{Z}^+$ denote how many times σ is emitted at the state q while the machine M_j emits w .
- Let $n_w(M_j, q) \in \mathbb{Z}^+$ denote how many times the state q is visited while the machine M_j emits w .

RELATIVE FREQUENCY OF EMISSION

- Let $m_w(M_j, q, \sigma) \in \mathbb{Z}^+$ denote how many times σ is emitted at the state q while the machine M_j emits w .
- Let $n_w(M_j, q) \in \mathbb{Z}^+$ denote how many times the state q is visited while the machine M_j emits w .

Then

$$\text{freq}_w(\sigma|M_j, q) = \frac{m_w(M_j, q, \sigma)}{n_w(M_j, q)}, \quad (3)$$

represents the relative frequency that M_j emits σ at q during emission of w .

RELATIVE FREQUENCY OF EMISSION

- Let $m_w(M_j, q, \sigma) \in \mathbb{Z}^+$ denote how many times σ is emitted at the state q while the machine M_j emits w .
- Let $n_w(M_j, q) \in \mathbb{Z}^+$ denote how many times the state q is visited while the machine M_j emits w .

Then

$$\text{freq}_w(\sigma | M_j, q) = \frac{m_w(M_j, q, \sigma)}{n_w(M_j, q)}, \quad (3)$$

represents the relative frequency that M_j emits σ at q during emission of w .

It is straightforward to lift this definition to data sequences $D = \langle w_1 \times, w_2 \times, \dots, w_{|D|} \times \rangle$ by letting $w = w_1 \times w_2 \times \dots w_{|D|} \times$.

EMPIRICAL MEAN OF CO-EMISSION PROBABILITIES

EMPIRICAL MEAN OF CO-EMISSION PROBABILITIES

$$\text{sumCoemit}_w(\sigma, M_j, q) = \sum_{i \text{ s.t. } q(j,i)=q} \text{Coemit}(\sigma, i).$$

EMPIRICAL MEAN OF CO-EMISSION PROBABILITIES

$$\text{sumCoemit}_w(\sigma, M_j, q) = \sum_{i \text{ s.t. } q(j,i)=q} \text{Coemit}(\sigma, i).$$

The *empirical mean of a co-emission probability* is defined as follows:

$$\overline{\text{Coemit}}_w(\sigma | M_j, q) = \frac{\text{sumCoemit}_w(\sigma, M_j, q)}{n_w(M_j, q)}, \quad (4)$$

EMPIRICAL MEAN OF CO-EMISSION PROBABILITIES

$$\text{sumCoemit}_w(\sigma, M_j, q) = \sum_{i \text{ s.t. } q(j,i)=q} \text{Coemit}(\sigma, i).$$

The *empirical mean of a co-emission probability* is defined as follows:

$$\overline{\text{Coemit}}_w(\sigma | M_j, q) = \frac{\text{sumCoemit}_w(\sigma, M_j, q)}{n_w(M_j, q)}, \quad (4)$$

This is the sample average of the co-emission probability when $q \in Q_j$ is visited.

MAIN THEOREM

MAIN THEOREM

Consider any parameter $T_j(q, \sigma)$ in PDFA M_j .

Theorem

$\partial P(D)/\partial T_j(q, \sigma) = 0$ holds for all j if and only if the following equation is satisfied for all $1 \leq j \leq K$:

$$\text{freq}_w(\sigma|M_j, q) = \overline{\text{Coemit}_w(\sigma|M_j, q)} .$$

EXAMPLE

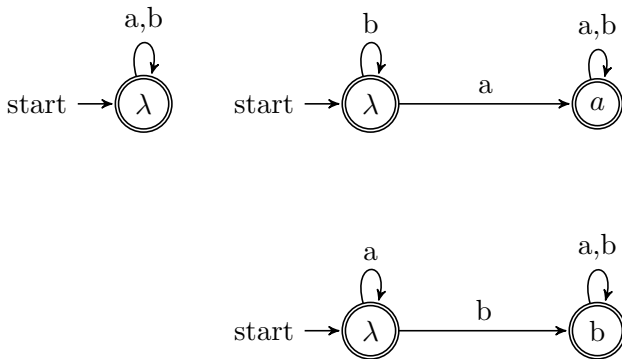


FIGURE: The 2-set of SD-PDFAs with $\Sigma = \{a, b\}$. There are 15 parameters. Suppose $D = abb \times bbb \times$.

EXAMPLE

$$\begin{array}{lll} \text{freq}_D(a|\mathcal{M}_\lambda, \lambda) = 1/8 & \text{freq}_D(a|\mathcal{M}_a, \lambda) = 1/5 & \text{freq}_D(a|\mathcal{M}_a, a) = 0/3, \\ \text{freq}_D(b|\mathcal{M}_\lambda, \lambda) = 5/8 & \text{freq}_D(b|\mathcal{M}_a, \lambda) = 3/5 & \text{freq}_D(b|\mathcal{M}_a, a) = 2/3, \\ \text{freq}_D(\times|\mathcal{M}_\lambda, \lambda) = 2/8 & \text{freq}_D(\times|\mathcal{M}_a, \lambda) = 1/5 & \text{freq}_D(\times|\mathcal{M}_a, a) = 1/3, \\ & \text{freq}_D(a|\mathcal{M}_b, \lambda) = 1/3 & \text{freq}_D(a|\mathcal{M}_b, b) = 3/5, \\ & \text{freq}_D(b|\mathcal{M}_b, \lambda) = 2/3 & \text{freq}_D(b|\mathcal{M}_b, b) = 0/5, \\ & \text{freq}_D(\times|\mathcal{M}_b, \lambda) = 0/3 & \text{freq}_D(\times|\mathcal{M}_b, b) = 2/5, \end{array}$$

FIGURE: Frequency computations with $D = abb \times bbb \times$ and the 2-set of SD-PDFAs on previous slide.

CONVEXITY OF THE NEGATIVE LOG LIKELIHOOD

CONVEXITY OF THE NEGATIVE LOG LIKELIHOOD

Let $\tau_{j,q,\sigma}$ denote $\log T_j(q, \sigma)$; i.e. the log of a parameter of C defined with $\bigotimes_j M_j$.

CONVEXITY OF THE NEGATIVE LOG LIKELIHOOD

Let $\tau_{j,q,\sigma}$ denote $\log T_j(q, \sigma)$; i.e. the log of a parameter of C defined with $\bigotimes_j M_j$.

Then τ can be thought of as a vector in \mathbb{R}^n where n is the number of parameters.

CONVEXITY OF THE NEGATIVE LOG LIKELIHOOD

Let $\tau_{j,q,\sigma}$ denote $\log T_j(q, \sigma)$; i.e. the log of a parameter of C defined with $\bigotimes_j M_j$.

Then τ can be thought of as a vector in \mathbb{R}^n where n is the number of parameters.

Theorem

$-\log P(w \times)$ is convex with respect to $\tau \in \mathbb{R}^n$.

CONVEXITY OF THE NEGATIVE LOG LIKELIHOOD

Let $\tau_{j,q,\sigma}$ denote $\log T_j(q, \sigma)$; i.e. the log of a parameter of C defined with $\bigotimes_j M_j$.

Then τ can be thought of as a vector in \mathbb{R}^n where n is the number of parameters.

Theorem

$-\log P(w_{\times})$ is convex with respect to $\tau \in \mathbb{R}^n$.

Thus the solution obtained by the previous theorem is a MLE.

DISCUSSION

DISCUSSION

At a high level, the problem we considered is a decomposition of complex probability distributions into simpler factors.

DISCUSSION

At a high level, the problem we considered is a decomposition of complex probability distributions into simpler factors.

This has also been studied in the context of Bayesian networks, Markov random fields, and probabilistic graphical models more generally (Bishop, 2006; Koller and Friedman, 2009).

DISCUSSION

At a high level, the problem we considered is a decomposition of complex probability distributions into simpler factors.

This has also been studied in the context of Bayesian networks, Markov random fields, and probabilistic graphical models more generally (Bishop, 2006; Koller and Friedman, 2009).

A reviewer points out that this literature may simplify our proofs.

FUTURE WORK

- ① Language Modeling with various sets of specific factors and various corpora such as ...
 - ① $SL_k + SP_k$
 - ② $SLP_{k,\ell}$ (Rogers and Lambert 2019, MoL)
 - ③ Atomic PDFAs based on phonological features (Chandlee et al. 2019, MoL)
- ② ... and compare to NNs, ALERGIA, and other algorithms on various benchmarks.
- ③ Connections to probabilistic graphical models
- ④ Extend results to weighted deterministic automata.

THANKS

We acknowledge Canaan Breiss, Morgan Cassels, Huteng Dai, Danny DeSantiago, Anton Kukhto, Jon Rawski, Yang Wang, and Yuhong Zhu for valuable feedback on a draft presentation.

Questions?