

Approximate Equivalence of Transduction

CMF 2007

ADRIEN VIEILLERIBIÈRE

Laboratoire de Recherche en Informatique
Université Paris Sud XI

May 21st, 2007

Motivations

- 1 Data Exchange
- 2 xslt Transformations
- 3 Model versus Specification

Outline: Preliminaries

Preliminaries

Equivalence

Trees

- 1 Word Edit Distance with Moves
- 2 Statistics
 - Uniform Statistic of a Word
 - Link Statistics-Distance
 - Uniform Statistic of a Loop
- 3 Word Transducer
 - Definition
 - Relation defined by a Transducer
 - Powering
 - Graph of the strongly connected components

Elementary operations

- Insert a letter : $aaabbb \xrightarrow{ins} aaa**c**bbb$
- Remove a letter : $aaacbbb**b** \xrightarrow{del} aaacbb$
- Change a label : $aaacbb**b** \xrightarrow{swi} aaacbc$
- Move a subword : $a**aac**bc \xrightarrow{mov} acb**aac**$

Edit Distance with moves

- 1 $\text{dist}(w_1, w_2) = \min\{p \mid w_1 \xrightarrow{o_1} \dots \xrightarrow{o_p} w_2\}$
- 2 $d(w_1, w_2) = \frac{\text{dist}(w_1, w_2)}{\max(|w_1|, |w_2|)}$
- 3 $\text{dist}(w_1, L) = \min_{w_2 \in L} \text{dist}(w_1, w_2)$
- 4 $\text{dist}(L_1, L_2) = \min_{w_1 \in L_1, w_2 \in L_2} \text{dist}(w_1, w_2)$

Edit Distance with moves

$$① \text{ dist}(w_1, w_2) = \min\{p \mid w_1 \xrightarrow{o_1} \dots \xrightarrow{o_p} w_2\}$$

$$② d(w_1, w_2) = \frac{\text{dist}(w_1, w_2)}{\max(|w_1|, |w_2|)}$$

$$③ \text{ dist}(w_1, L) = \min_{w_2 \in L} \text{dist}(w_1, w_2)$$

$$④ \text{ dist}(L_1, L_2) = \min_{w_1 \in L_1, w_2 \in L_2} \text{dist}(w_1, w_2)$$

Examples

$$① \text{ dist}(aaabbb, bbbaac) = 2$$

$$② d(aaabbb, bbbaac) = 1/3$$

$$③ \text{ dist}(aaabbb, a^*) = 3$$

$$④ \text{ dist}((aaa)^+.(bbbbbbb)^+, (bbbbbbb)^*.(aaa)^+) = 1$$

1 Word Edit Distance with Moves

2 Statistics

- Uniform Statistic of a Word
- Link Statistics-Distance
- Uniform Statistic of a Loop

3 Word Transducer

- Definition
- Relation defined by a Transducer
- Powering
- Graph of the strongly connected components

Let Σ be a binary alphabet, $w \in \Sigma^*$ of length n and $k \in \mathbb{N}^+$

Statistic of a Word : ustat_k

$$\text{ustat}_k(w) = \frac{1}{n - k + 1} \cdot \begin{pmatrix} \#n_1 \\ \#n_2 \\ \vdots \\ \#n_{2^k} \end{pmatrix} \text{ where}$$

$\#n_1$: number of $00\dots 0$
 $\#n_2$: number of $0\dots 01$
 \vdots
 $\#n_{2^k}$: number of $1\dots 11$

Example

$w = 001010101110$, length $n = 12$, $k = 2$, $n - k + 1 = 11$ blocks

$$\text{ustat}(w) = \frac{1}{11} \begin{pmatrix} 1 \\ 4 \\ 4 \\ 2 \end{pmatrix}$$

Let Σ be a alphabet, $w \in \Sigma^*$ of length n and $k \in \mathbb{N}^+$

Statistic of a Word : $ustat_k$

$$ustat_k(w) = \frac{1}{n - k + 1} \cdot \begin{pmatrix} \#n_1 \\ \#n_2 \\ \vdots \\ \#n_{|\Sigma|^k} \end{pmatrix}$$

Main Properties

- 1 $d(w, w') \approx |\text{ustat}(w) - \text{ustat}(w')|_1$ for words of similar length
- 2 $|\text{ustat}(w) - Y(w)|_1 \leq \varepsilon$ for $Y(W)$ the ustat vector on N samples
- 3 distance between words
 - NP-complete
 - Testable, $\mathcal{O}(1)$:
Sample N subwords of length k : $Y(w)$ and $Y(w')$.
If $|Y(w) - Y(w')|_1 < \varepsilon = \frac{1}{k}$ accept, else reject

Limit Statistic of a loop

Let $L = w^*$, $|w| = q$

$$\mathcal{H}_{\text{lim}}(L)[u] = \Pr_{j=1, \dots, q} \left[w[j].w[j+1 \text{ Mod } k] \dots w[j+k-1 \text{ Mod } k] = u \right]$$

Example

$$L = (0001)^*, k = 2, \mathcal{H}_{\text{lim}}(L) = \begin{pmatrix} 1/2 \\ 1/4 \\ 1/4 \\ 0 \end{pmatrix}$$

- 1 Word Edit Distance with Moves
- 2 Statistics
 - Uniform Statistic of a Word
 - Link Statistics-Distance
 - Uniform Statistic of a Loop
- 3 **Word Transducer**
 - **Definition**
 - **Relation defined by a Transducer**
 - **Powering**
 - **Graph of the strongly connected components**

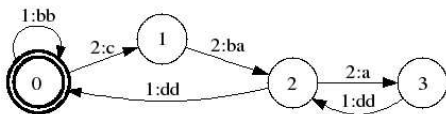
Word Transducers

A word transducer \mathcal{T} is a tuple $(\Sigma, \Sigma', Q, I, F, \delta)$ where

- Σ in the input alphabet, Σ' the output alphabet
- Q a set of states, $I \in Q$ an initial state, $F \subseteq Q$ final states
- $\delta : Q \times (\Sigma \cup \Lambda) \times (\Sigma')^* \times Q$ the transition fonction

Example of a Word Transducer

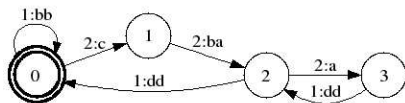
- $\Sigma = \{1, 2\}, \Sigma' = \{a, b, c, d\}$
- $Q = \{0, 1, 2, 3\}, I = \{0\}, F = \{0\}$
- $\delta = \{(0, 1, bb, 0), (0, 2, c, 1), (1, 2, ba, 2), (2, 1, dd, 0), (2, 2, a, 3), (3, 1, dd, 2)\}$



Relation of a Word Transducer \mathcal{T}

For two words w_1 and w_2 , (w_1, w_2) is in the relation defined by \mathcal{T} , if there is a run from the initial state to a final state inputing w_1 and outputing w_2 .¹

Example

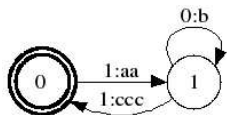


(Λ, Λ) , $(1, bb)$, $(22211, cbaadddd)$ are in the relation defined by \mathcal{T}

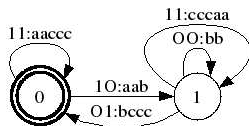
¹We then note $(w_1, w_2) \in \mathcal{T}$

Example : Squaring

a transducer



its square



k -Powering of a Transducer

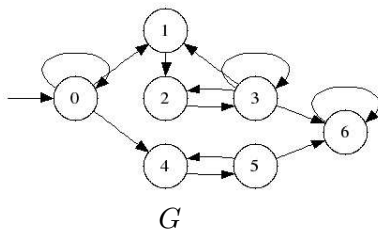
Let $\mathcal{T} = (\Sigma, \Sigma', Q, I, F, \delta)$ a transducer and $k \in \mathbb{N}$.

$\mathcal{T}^k = (\Sigma^k, \Sigma', Q, I, F, \delta^k)$ where

$(q_1, a_1 a_2 \dots a_k, w, q_{k+1}) \in \delta^k$ iff there exist $q_1, \dots, q_k \in Q$ and $w_1, \dots, w_k \in (\Sigma')^*$ such that $w = w_1.w_2\dots w_k$ and for all i in $\{1, \dots, k\}$, $(q_i, a_i, w_i, q_{i+1}) \in \delta$.

For G the graph of a (delabeled) automaton or transducer, \hat{G} denotes the graph of strongly connected components of G

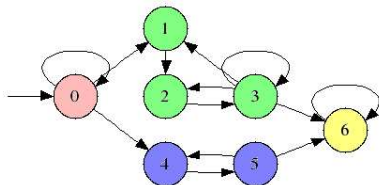
Example



For G the graph of a (delabeled) automaton or transducer, \hat{G} denotes the graph of strongly connected components of G

Example

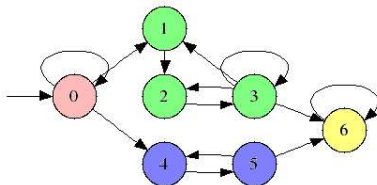
G



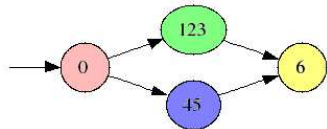
For G the graph of a (delabeled) automaton or transducer, \hat{G} denotes the graph of strongly connected components of G

Example

G



\hat{G}



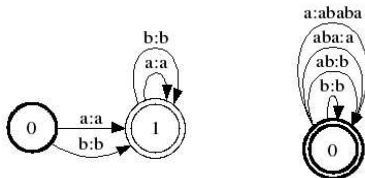
Outline: Equivalences

- 4 warm up
- 5 definition
- 6 exact case
- 7 ε -equivalence
 - two transducers ε -close
 - statistical representation
 - In the relation \Rightarrow close to the embedding
 - In the embedding \Rightarrow close to the relation
 - counterexamples : sources of non ε -equivalence

The intersection is undecidable

Reduction to the Post Correspondence Problem

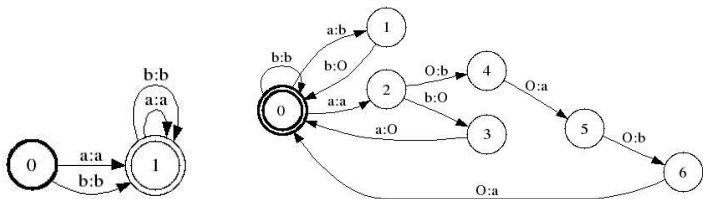
- an instance: $\left(\begin{array}{c} a \\ ababa \end{array} \right), \left(\begin{array}{c} aba \\ a \end{array} \right), \left(\begin{array}{c} ab \\ b \end{array} \right), \left(\begin{array}{c} b \\ b \end{array} \right)$
- a solution: $\left(\begin{array}{c} aba \\ a \end{array} \right) \left(\begin{array}{c} b \\ b \end{array} \right) \left(\begin{array}{c} a \\ ababa \end{array} \right) \left(\begin{array}{c} b \\ b \end{array} \right) \left(\begin{array}{c} aba \\ a \end{array} \right)$



The intersection is undecidable

Reduction to the Post Correspondence Problem

- an instance: $\left(\begin{array}{c} a \\ ababa \end{array} \right), \left(\begin{array}{c} aba \\ a \end{array} \right), \left(\begin{array}{c} ab \\ b \end{array} \right), \left(\begin{array}{c} b \\ b \end{array} \right)$
- a solution: $\left(\begin{array}{c} aba \\ a \end{array} \right) \left(\begin{array}{c} b \\ b \end{array} \right) \left(\begin{array}{c} a \\ ababa \end{array} \right) \left(\begin{array}{c} b \\ b \end{array} \right) \left(\begin{array}{c} aba \\ a \end{array} \right)$



From $\mathcal{T}_1 \subseteq \mathcal{T}_2$ to $\mathcal{T}_1 \subseteq_\varepsilon \mathcal{T}_2$

- 1 **exact inclusion:** $(w, w') \in \mathcal{T}_1 \implies (w, w') \in \mathcal{T}_2$
- 2 **ε -inclusion:** $(w, w') \in \mathcal{T}_1 \implies \exists (\tilde{w}, \tilde{w}') \in \mathcal{T}_2$ such that, for the classical edit distance with moves
 - \tilde{w} is ε -close to w i.e. $d(w, \tilde{w}) \leq \varepsilon$
 - \tilde{w}' is ε -close to w' i.e. $d(w', \tilde{w}') \leq \varepsilon$

From $\mathcal{T}_1 \subseteq \mathcal{T}_2$ to $\mathcal{T}_1 \subseteq_\varepsilon \mathcal{T}_2$

- 1 **exact inclusion:** $(w, w') \in \mathcal{T}_1 \implies (w, w') \in \mathcal{T}_2$
- 2 **ε -inclusion:** $(w, w') \in \mathcal{T}_1 \implies \exists (\tilde{w}, \tilde{w}') \in \mathcal{T}_2$ such that, for the classical edit distance with moves
 - \tilde{w} is ε -close to w i.e. $d(w, \tilde{w}) \leq \varepsilon$
 - \tilde{w}' is ε -close to w' i.e. $d(w', \tilde{w}') \leq \varepsilon$

Finer ? Coarser ?

exact inclusion $\implies \varepsilon$ -inclusion

Equivalence is undecidable

- 1 The equivalence problem for non-deterministic transducers (one way, finite state) is undecidable even when they cannot read or write the empty string

T. V. Griffiths, "The unsolvability of the Equivalence Problem for Lambda-Free nondeterministic generalized machines", 1968.

- 2 The equivalence problem for 1-free² input deterministic³ sequential transducers is undecidable

J. Karhumäki and L. P. Lisovik, "On the Equivalence of Finite Substitutions and Transducers", 1999, (Corollary 1).

²Transitions do not output the empty string

³The (input) underlying automaton is deterministic

Decidability

- 1 The equivalence for deterministic transducers is decidable
Bird 1973, Vaillant 1974
- 2 The equivalence for n-tape deterministic finite automata is decidable
Harju & Karhumäki, 1991
- 3 The equivalence for sequential transducers satisfying the prefix condition⁴ is decidable
Karhumäki & Lisovik, 1999
- 4 The equivalence on an NR set of strings for two deterministic two-way finite state transducers is decidable
Engelfriet, Maneth, 2005

⁴Avoid two transitions from the same state, the same input letter, and an output different and prefix of the other

PSPACE-hardness

Equivalence is at least PSPACE-hard, even restricted to identity edges

proof

as hard as the equivalence of non-deterministic finite automata.

PSPACE-completeness for a weak form of non-determinism

The equivalence of deterministic two-way (sequential) transducers which are allowed to make some finite number of non-deterministic moves is a PSPACE-complete problem pointed out by Gurari, 1982

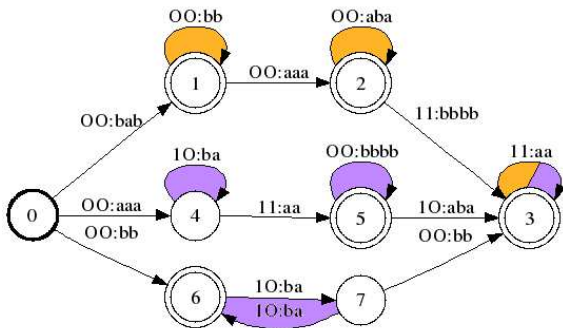
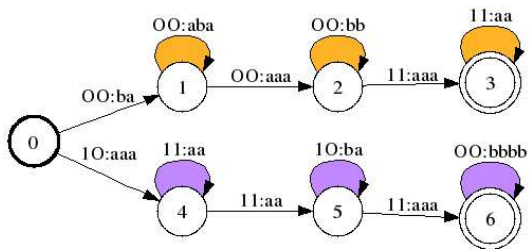
4 warm up

5 definition

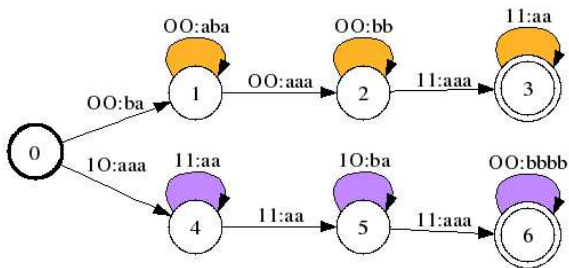
6 exact case

7 ϵ -equivalence

- two transducers ϵ -close
- statistical representation
- In the relation \Rightarrow close to the embedding
- In the embedding \Rightarrow close to the relation
- counterexamples : sources of non ϵ -equivalence

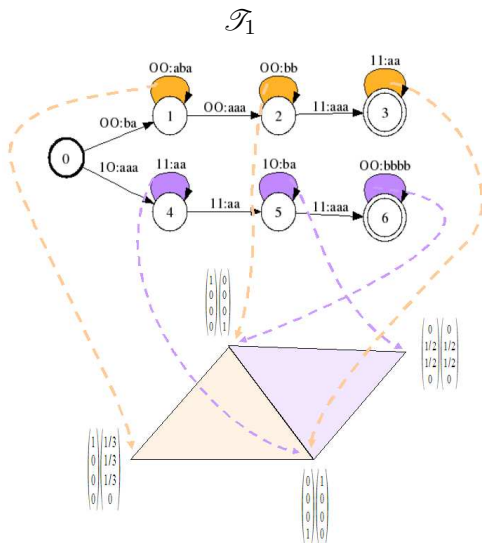


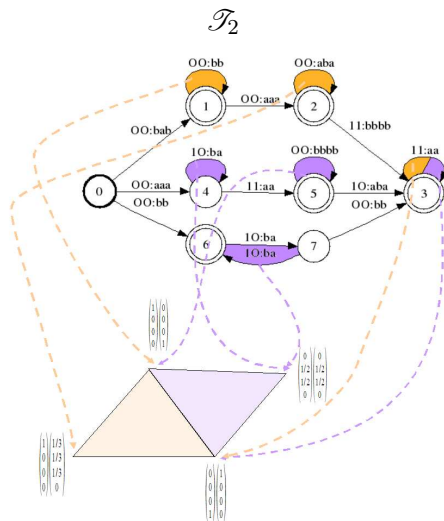
Forget finite contributions



In essence:

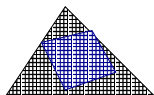
$$t = (00 : aba)^* . (00 : bb)^* . (11 : aa)^* + (11 : aa)^* . (10 : ba)^* . (00 : bbbb)^*$$





ε -Inclusion

- Apply the k -powering of \mathcal{T}_1 and \mathcal{T}_2 , $k = 1/\varepsilon$:
 - Matrix multiplication
- compute the embeddings $\mathcal{H}_\varepsilon = (H_\varepsilon^I, H_\varepsilon^O)$ in the statistical space:
 - analyze paths π in $\widehat{\mathcal{T}}_i^k$
 - convex hulls for π -compatible loops
- $H_\varepsilon^1 \subseteq H_\varepsilon^2$?:
 - H_ε^i is a union of polytopes
 - inclusion of the grids of step ε



- check lengths

\implies PSPACE tester for every fixed $\varepsilon > 0$

Assume $(w, w') \in \mathcal{T}^k$

- let π be a path in $\widehat{\mathcal{T}}^k$ s.t. $w' \in \mathcal{T}_\pi^k(w)$
- following π , $w = v_0 u_1 u_2 \dots u_l v_l$,
 $|v_0|, |u_1|, \dots, |u_l|, |v_l| \leq m.k$, $\{u_1, \dots, u_l\}$: π -compatible
- $\{u_1, \dots, u_l\} = \{u'_1, \dots, u'_b\}$ with $b \leq a^k$
- take $\tilde{w} = v_0 (u'_1)^{k_1} v_1 (u'_2)^{k_2} \dots v_{b-1} (u'_b)^{k_b} v_b$ with
 $v_0 v_1 \dots v_b$: witness from compatibility,
 v_i : sequences of less than m edges.
- $\text{dist}(w, \tilde{w}) \leq (b + 1)m + \varepsilon n$
- take $\tilde{w}' = \text{out}(v_0) (\text{out}(u'_1))^{k_1} \text{out}(v_1) \dots (\text{out}(u'_b))^{k_b} \text{out}(v_b)$
- $(\tilde{w}, \tilde{w}') \in \mathcal{T}^k$
- $\text{dist}(w', \tilde{w}') \leq \max_{\text{out}}(b + 1)m + \varepsilon n$

$(\text{ustat}(\tilde{w}), \text{ustat}(\tilde{w}'))$ is close to the polytope

- $\tilde{x} = \text{ustat}(\tilde{w}), \tilde{x}' = \text{ustat}(\tilde{w}')$
- $\beta_i = \frac{k_i}{\sum_i k_i}$
- $\beta = (\beta_1, \dots, \beta_b)$
- $x = \text{Renorm}(\sum_{i=1}^b \beta_i l_{in}(B_i) \mathcal{H}_{lim_{in}}(B_i))$
- $x \in \mathcal{H}^I$
- $|x - \tilde{x}| \leq \frac{(b+1)(3(k-1)+m.p)}{\tilde{n}}$
- $x' = \text{Renorm}(\sum_{i \in b} \frac{l_{out}(B_i)}{l_{in}(B_i)} \beta_i \mathcal{H}_{lim_{out}}(B_i))$
- $(x, x') \in \mathcal{H}$
- $|x' - \tilde{x}'| \leq \frac{\text{maxout}((b+1)(3(k-1)+m.p))}{\tilde{n}'}$

Assume $(\lambda, \lambda') \in \mathcal{H}$

- take sets of summits T , (λ, λ', π) -compatibles:

- $\lambda = \sum_{i \in T} \lambda_i \mathcal{H}_{lim_{in}}(B_i)$

- $\lambda' = \sum_{i \in T} \lambda'_i \mathcal{H}_{lim_{out}}(B_i)$

- $\lambda' = Renorm \left(\frac{l_{out}(B_i)}{l_{in}(B_i)} \lambda_i \right)$

- approximate λ_i by rationals α_i

- denormalize $\tilde{\lambda}_i = c \cdot \frac{l_{out}(B_i)}{l_{in}(B_i)} \alpha_i$

- take $w = u_0(in(B_1))^{\tilde{\lambda}_1} u_1(in(B_2))^{\tilde{\lambda}_2} u_2 \dots u_{|T|}$

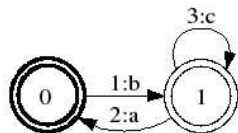
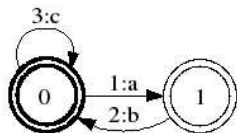
- take $w' = u'_0(out(B_1))^{\tilde{\lambda}_1} u'_1(out(B_2))^{\tilde{\lambda}_2} u'_2 \dots u'_{|T|}$

- $(w, w') \in \mathcal{T}$

- $|(ustat(w), ustat(w')) - (\lambda, \lambda')|_1 \leq 2(maxout + 1)\varepsilon$

Sources of non ϵ -equivalence

1 statistical domains ϵ -far

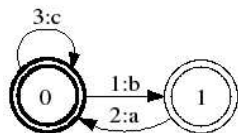
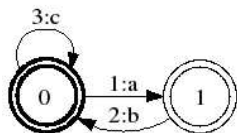


2 statistical images ϵ -far

3 generated length asymptotically different

Sources of non ϵ -equivalence

- 1 statistical domains ϵ -far
- 2 statistical images ϵ -far



- 3 generated length asymptotically different

Sources of non ϵ -equivalence

- 1 statistical domains ϵ -far
- 2 statistical images ϵ -far
- 3 generated length asymptotically different



Sources of non ε -equivalence

- 1 statistical domains ε -far
- 2 statistical images ε -far
- 3 generated length asymptotically different

Solutions :

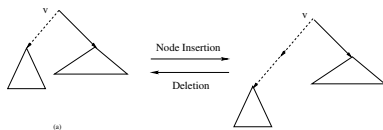
- 1 the domain is regular: similar embedding of a regular language: PTIME tester.
- 2 use the embedding
- 3 generating functions (rational functions)

- 8 Distances and Statistics
 - Tree Edit Distance with Moves

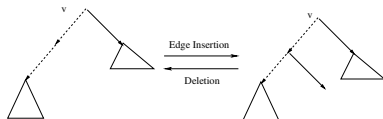
- 9 Tree Transducers
 - Definition
 - Example
 - the next step

Tree Edit Distance with Moves : Elementary Operations

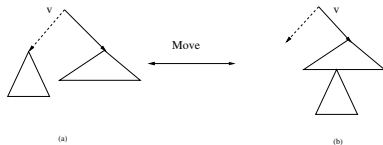
Node operation:



Edge operation:



Move:



Relabeling

Let $H_{\Sigma'}[Q]$ be the set of finite sequences of finite trees where one leaf of each tree is a distinguished element labeled by a sequence of states in Q , which is possibly empty.

Tree Transducers

A tree transducer \mathcal{T} is a tuple $(\Sigma, \Sigma', Q, I, \delta)$ where

- Σ in the input alphabet, Σ' the output alphabet
- Q a set of states, $I \in Q$ an initial state
- $\delta : \Sigma \times Q \rightarrow H_{\Sigma'}[Q]$,

An xslt program

```
<xsl:output method="xml" indent="yes"/>

<xsl:template match="/">
  <root> <s> <summary/>
    <xsl:apply-templates mode="q0"/>
  </s>
  <c> <xsl:apply-templates mode="q1"/>
  </c>
</root>
</xsl:template>

</xsl:template match="a" mode="q0">
  <a/>
</xsl:template>

</xsl:template match="b" mode="q0">
  <b/>
</xsl:template>

</xsl:template match="a" mode="q1">
  <a>
    <xsl:apply-templates mode="q1"/>
  </a>
</xsl:template>

</xsl:template match="b" mode="q1">
  <b>
    <xsl:apply-templates mode="q1"/>
  </b>
</xsl:template>
```

- 1 define tree statistics as matrices
- 2 Matrices distances & Edit distance
- 3 applications for xslt programs

thank you

Questions ?

vieille@lri.fr