

Edit distance of finite state transducers

Saina Sunny



PhD Scholar, IIT Goa

Joint work with **C. Aiswarya** (CMI) and **Amaldev Manuel** (IIT Goa)

Transducers

• Transducers are finite state automata that produce output words



Delete all *a*'s before first *b*

 $aaba \rightarrow ba$

 $bab \rightarrow bab$

Transducers

• Transducers are finite state automata that produce output words





Delete all *a*'s before first *b*

How do we compare transducers?

- Checking equivalence of two transducers
 - decidable for rational functions [Gurari-Ibarra' 1983],
 - decidable for regular functions [Gurari'1982,Culik-Karhumaki'1987]
 - open for polyregular functions [Bojanczyk'2018]
 - undecidable for rational relations [Griffiths'1968]
- Can we say something meaningful about non-equivalent transducers?

How do we compare transducers?



• Functional equivalence (on any input, the respective outputs are "exactly" the same)

How do we compare transducers?

• Relax it : on any input, the respective outputs are close enough





Metric on transducers

• Let d be a metric on words. Lift it to word-to-word functions (transductions).

$$d(T_1, T_2) = \begin{cases} \sup \{d(T_1(w), T_2(w))\} \\ \infty \end{cases}$$

- T_1 and T_2 are close if $d(T_1, T_2)$ is finite.



 $(w)) \mid w \in dom(T_1)\} \quad \text{if } dom(T_1) = dom(T_2)$ otherwise

• Related work: adjacent functions [Reutenauer-Schützenberger '1991]

Edit Distances

- Given a set of edit operations,

• Edit distances between two words is the minimum number of edits required to convert one to another.

> ababa hahah

• Ex: insert a letter, delete a letter, or substitute a letter with another



Common Edit distances

Edit Distances
Hamming distance
Transposition distance
Conjugacy distance
Levenshtein edit distance
Longest common subsequence
Damerau-Levenshtein distance

Edit operations
 letter-to-letter substitution
swapping adjacent letters
left and right cyclic shifts
insertion, deletion, substitution
insertion and deletion
Insertion, deletion, substitution and adjacent transposition

Edit distances - preorder relation

Edit Distances	Edit opera
Hamming distance	letter-to-letter
Transposition distance	swapping adja
Conjugacy distance	left and right c
Levenshtein edit distance	insertion, d substitu
Longest common subsequence	insertion and
Damerau-Levenshtein distance	Insertion,deletion and adjacent tr



Common Edit distances

Edit Distances
Hamming distance
Transposition distance
Conjugacy distance
Levenshtein edit distance
Longest common subsequence
Damerau-Levenshtein distance
Discrete

Edit operations
letter-to-letter substitution
swapping adjacent letters
 left and right cyclic shifts
insertion, deletion, substitution
insertion and deletion
Insertion, deletion, substitution and adjacent transposition
 Ø

Metric on transducers Example



- For each block of *a*, output a
- For each block of *b*, output b

•
$$d_{lev}(T_1, T_2) = 2$$



• For each block of *a*, output b

• For each block of *b*, output a

 $aaabbabbba \rightarrow (ababa, babab)$

- $d(T_1, T_2) = \infty$ if only
- substitutions
- cyclic shifts
- adjacent swapping



Metric on transducers Questions

$$d(T_1, T_2) = \begin{cases} \sup \{d(T_1(w), T_2(w)) \mid 1 \\ \infty \end{cases}$$

- Given T_1 , T_2 is $d(T_1, T_2)$ computable?
- Given T_1 , T_2 is $d(T_1, T_2)$ finite?
- Given T_1, T_2 and $k \in \mathbb{N}$, is $d(T_1, T_2)$ at most k?

 $w \in dom(T_1)$ if $dom(T_1) = dom(T_2)$ otherwise

(Distance)

(Closeness)

(k-closeness)

Metric on transducers Results

Problem	Input	Question
Distance Problem	transducers $\mathcal{T}_1, \mathcal{T}_2$	$d(\mathcal{T}_1,\mathcal{T}_2)?$
Closeness Problem	transducers $\mathcal{T}_1, \mathcal{T}_2$	Is $d(\mathcal{T}_1, \mathcal{T}_2) < \infty$?
k-closeness Problem	integer k, transducers $\mathcal{T}_1, \mathcal{T}_2$	Is $d(\mathcal{T}_1, \mathcal{T}_2) \leq k$?

<u>Proposition</u>: Distance is computable iff closeness and k-closeness is decidable for integer-valued metrics

<u>Theorem</u>: Closeness and k-closeness for rational functions is decidable for all metrics $d \in \{d_{lev}, d_{lcs}, d_{damerau}, d_{conj}, d_{ham}, d_{trans}\}$.

Closeness and k-closeness

- Given transducers T_1, T_2
 - Domain of T_1 and T_2 must be same.
 - Let *T* be the cartesian product of T_1 and T_2

Cartesian product of two transducers



Output *a*'s before *b*



Output *a*'s after *b*



 $aaba \rightarrow (a, \epsilon) \cdot (a, \epsilon) \cdot (\epsilon, \epsilon) \cdot (\epsilon, a) = (aa, a)$



Closeness and k-closeness

- Given transducers T_1, T_2
 - Domain of T_1 and T_2 must be same.
 - Let *T* be the cartesian product of T_1 and T_2
 - generates set of all pairs of output words of T_1 , T_2 on any input
 - Loops of T must generate output pairs of same length (Close w.r.t. d_{len})

k-closeness **For edit distances**

- Given transducers T_1, T_2
 - Domain of T_1 and T_2 must be same.
 - Let *T* be the cartesian product of T_1 and T_2 — generates set of all pairs of output words of T_1, T_2 on any input
 - Loops of T must generate output pairs of same length (Close w.r.t. d_{lon})
 - From T, construct an automaton that accepts w if $d(T_1(w), T_2(w)) \leq k$ 1.
 - 2. Start with budget k. Non-deterministically do edits, update the budget and residues appropriately. Budget is not allowed to be negative.
 - 3. Check if the language accepted is the domain of T. Yes: k-close; No: not k-close.

Closeness

- Given transducers T_1, T_2
 - Domain of T_1 and T_2 must be same.
 - Let *T* be the cartesian product of T_1 and T_2
 - generates set of all pairs of output words of T_1 , T_2 on any input
 - Loops of *T* must generate output pairs of same length (Close w.r.t. d_{len})
 - Identical output pairs?





if u = xy and v = yx

 $\begin{array}{c} x | y x y \cdots x y x y \\ y x y x \cdots y x y | x \end{array}$

• *u* and *v* are conjugates if u = xy and v = yx





 $\begin{array}{c} x \mid y \mid x \mid y \cdots \mid x \mid y \mid x \mid y \\ y \mid x \mid y \mid x \cdots \mid y \mid x \mid y \mid x \\ \end{array}$

- *u* and *v* are conjugates if u = xy and v = yx
- Equivalently, if there exists a word z such that





X Yy X



- *u* and *v* are conjugates if u = xy and v = yx
- Equivalently, if there exists a word z such that



y and v = yxuch that uz = zv





- *u* and *v* are conjugates if u = xy and v = yx
- Equivalently, if there exists a word *z* such that





 $\mathcal{U} \mathcal{Z} = \mathcal{Z} \mathcal{V}$

- Given transducers T_1, T_2
 - Domain of T_1 and T_2 must be same.
 - Let T be the cartesian product of T_1 and T_2
 - generates set of all pairs of output words of T_1 , T_2 on any input
 - Loops of T must generate output pairs of same length

<u>Lemma</u>: If T_1 and T_2 are close w.r.t metric $d \in \{d_{lev}, d_{lcs}, d_{damerau}, d_{conj}, d_{ham}, d_{trans}\}$, then every loop in T generates only conjugate pair of words.





More on conjugates



(*abca*, *baac*) is not conjugate



Any combination of pairs is conjugate

 $a/b \cdots a b/a/c \cdots a c/b/a/c \cdots$ $pa\cdots p|a| ca\cdots c|a|$

u z = z v

More on conjugates

- G set of pairs of words
- G^* consist of pairs obtained by point wise concatenation of some pairs in G

When is every pair in G^* conjugate?

Theorem[Aiswarya-Manuel-S. '2024] The set G^* is conjugate iff G has a common witness





Closeness **For Levenshtein family**

- Given transducers T_1, T_2
 - Domain of T_1 and T_2 must be same.
 - Let T be the cartesian product of T_1 and T_2

Lemma: T_1 and T_2 are close w.r.t. Levenshtein distance if and only if all the loops of T generate only conjugate words



— generates set of all pairs of output words of T_1, T_2 on any input

More on conjugates

Theorem[Aiswarya-Manuel-S. '2024] The set $(u_0, v_0)G_1^*(u_2, v_2)G_2^*\cdots G_k^*(u_k, v_k)$ is conjugate iff it has a common witness



Closeness **For Conjugacy**

- Given transducers T_1, T_2
 - Domain of T_1 and T_2 must be same.
 - Let T be the cartesian product of T_1 and T_2

<u>Proposition</u>: T_1 and T_2 are close w.r.t. conjugacy distance if and only if T generate only conjugate words

• Conjugacy of a rational relation is decidable [Aiswarya-Manuel-S. '2024]

– generates set of all pairs of output words of T_1, T_2 on any input



Closeness For Hamming and transposition

- Output lengths must be equal for all words.
 - Loops generate words of same length
 - Delay between partial outputs depends only on the state.
- Check if all pairs in a loop (modulo the delay on border) are equal (up to some length depending on the delay).



Metric on transducers Results

Problem	Input	Question
Distance Problem	transducers $\mathcal{T}_1, \mathcal{T}_2$	$d(\mathcal{T}_1,\mathcal{T}_2)?$
Closeness Problem	transducers $\mathcal{T}_1, \mathcal{T}_2$	Is $d(\mathcal{T}_1, \mathcal{T}_2) < \infty$?
k-closeness Problem	integer k, transducers $\mathcal{T}_1, \mathcal{T}_2$	Is $d(\mathcal{T}_1, \mathcal{T}_2) \leq k$?

<u>Proposition</u>: Distance is computable iff closeness and k-closeness is decidable for integer-valued metrics

<u>Theorem</u>: Closeness and k-closeness for rational function is decidable for all metrics $d \in \{d_{lev}, d_{lcs}, d_{damerau}, d_{conj}, d_{ham}, d_{trans}\}$.

Related notions and generalisations

Diameter of a Rational Relation

• The diameter of a rational relation R w.r.t. a metric d is the supremum of distance of each pair of words in R

$$dia_d(R) = \sup\{d(u, v) \mid (u, v) \in R\}$$

• Related Work: rational relation with bounded delay [Frougny-Sakarovitch'1991]

Diameter of a Rational Relation Questions

Problem	Input	Question
Diameter Problem	rational relation R	$dia_d(R)?$
Bounded Diameter Problem	rational relation R	Is $dia_d(R) < \infty$?
$k\mathchar`-bounded$ Diameter Problem	integer k , rational relation R	Is $dia_d(R) \le k?$



Diameter of a Rational Relation

Problem	Input	Question
Diameter Problem	rational relation ${\cal R}$	$dia_d(R)?$
Bounded Diameter Problem	rational relation ${\cal R}$	Is $dia_d(R) < \infty$?
$k\mathchar`-bounded$ Diameter Problem	integer k , rational relation R	Is $dia_d(R) \le k?$

<u>Proposition</u>: Diameter problem of a rational relation is mutually reducible to distance problem of two rational functions

Results



<u>Proposition</u>: Diameter problem of a rational relation is mutually reducible to distance problem of two rational functions

- Distance -> Diameter
 - Given two transducers T_1, T_2 , check if their domains are equal
- Diameter -> Distance
 - By virtue of [Nivat'1968] theorem

• $d(T_1, T_2) = dia_d(R)$ where R is the relation generated by cartesian product of T_1 and T_2

Diameter of a Rational Relation

Problem	Input	Question
Diameter Problem	rational relation R	$dia_d(R)?$
Bounded Diameter Problem	rational relation ${\cal R}$	Is $dia_d(R) < \infty$?
$k\mathchar`-bounded$ Diameter Problem	integer k , rational relation R	Is $dia_d(R) \le k?$

<u>Proposition</u>: Diameter problem of a rational relation is mutually reducible to distance problem of two rational functions

<u>Corollary</u>: All the above problems are decidable for rational relation w.r.t. metrics $d \in \{d_{lev}, d_{lcs}, d_{damerau}, d_{conj}, d_{ham}, d_{trans}\}$

Results



Index of relation in a composition closure

• Index of a rational relation R in the composition closure of S is the smallest integer k such that R is contained in at most k-fold composition of S

$$R \subseteq \bigcup_{0 \le i \le k} \underbrace{S \circ S}_{i \text{ times}}$$

• Example:

 $\{a,b\}^* \times \{a,b\}^*$

- S deletes the first a if exists on any input
- R_k deletes first k a's if exist on any input
- R delete all *a*'s on any input

••• o S

nes

- $Index(R_k, S) = k$
- $Index(R, S) = \infty$

Index of relation in a composition closure Questions

Problem	Input	Question
Index Problem	rational relation R, S	$\operatorname{Index}(R,S)?$
Bounded (or Finite) Index Problem	rational relation R, S	Is $Index(R, S) < \infty$?
$k\mbox{-}{\rm bounded}$ Index Problem	integer k , rational relation R , S	Is $Index(R, S) \le k?$

Index of relation in a composition closure Results

ProblemInputIndex ProblemrationalBounded (or Finite) Index Problemrationalk-bounded Index Probleminteger

<u>Lemma</u>: It is undecidable to check if a rational relation has a bounded index in the composition closure of an arbitrary rational relation

	Question
al relation R, S	$\operatorname{Index}(R,S)?$
al relation R, S	Is $Index(R, S) < \infty$?
r k , rational relation R , S	Is $Index(R, S) \le k$?

Metrizable Relation

- Graph of a relation S vertices (words), edge (between related words in S)
- $d_{S}(u, v) =$ length of the shortest path between u and v in the graph of S
- S is d-metrizable if d_S is equivalent to metric d up to boundedness.

Proposition: The index of a rational relation in the composition closure of a dmetrizable relation is computable for $d \in \{d_{len}, d_{lev}, d_{lcs}, d_{dl}, d_h, d_{trans}, d_{coni}\}$

Index of relation in a composition closure **Results**

Problem	Input	Question
Index Problem	rational relation R, S	$\operatorname{Index}(R,S)?$
Bounded (or Finite) Index Problem	rational relation R, S	Is $Index(R, S) < \infty$?
$k\mbox{-}{\rm bounded}$ Index Problem	integer k , rational relation R , S	Is $Index(R, S) \le k$?

Lemma: It is undecidable to check if a rational relation has a bounded index in the composition closure of an arbitrary rational relation

<u>Corollary</u>: All the above problems are decidable for rational relation in the composition closure of d- metrizable relation for $d \in \{d_{len}, d_{lev}, d_{lcs}, d_{dl}, d_{h}, d_{trans}, d_{conj}\}$

Conclusion

- We have defined the following notions
 - Distance between rational functions
 - Diameter of rational relation
 - Index of a rational relation in a composition closure

• All are computable w.r.t. metrics $d \in \{d_{len}, d_{lev}, d_{lcs}, d_{dl}, d_h, d_{trans}, d_{coni}\}$

Thank you

References

- 1. C. Aiswarya, Amaldev Manuel, Saina Sunny. Deciding conjugacy of a rational relation. CoRR, abs/ 2307.06777, 2023 (DLT to appear)
- 2. Christiane Frougny and Jacques Sakarovitch. Rational relations with bounded delay. In STACS 1991, volume 480 of LNCS, pages 50-63, 1991.
- 3. Christophe Reutenauer and Marcel-Paul Shutzenberger. Minimisation of rational word functions. SIAM Journal on Computing, 20(4):669-685, 1991.
- 4. Eitan M.Gurari. The equivalence problem for deterministic two-way sequential transducers is decidable. SIAM Journal on Computing, 11(3):448-452, 1982
- 5. Eitan M. Gurari and Oscar H.Ibarra. A note on finitely-valued and finitely ambiguous transducers. Math.Syst.Theory 16(1):61-66, 1983



References

- 6.
- Mikolaj Bojanczyk. Polyregular functions. CoRR abs/1810.08760, 2018. 7.
- Karel Culik, and Juhani Karhumaki. The equivalence problem for single-valued two-way transducers (on 8. NPDTOL Languages) is decidable. SIAM Journal on Computing, 1987.
- 9. Math. J, 9(4):289-298, 1962.
- generalized machines. J. ACM 15(3):409-413,1968

Maurice Nivat. Transduction des langages de Chomsky. PhD Thesis, Annales de l'Institut Fourier, 1968.

Roger C Lyndon, Marcel-Paul Schützenberger, et al. The equation $a^M = b^N c^P$ in a free group. Michigan

10. Timothy V.Griffiths. The unsolvability of the equivalence problem for lambda-free nondeterministic