

# SEDiL: Software for Edit Distance Learning

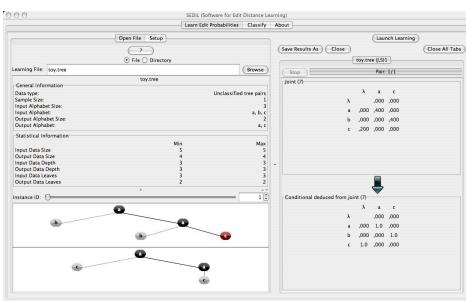
<http://labh-curien.univ-st-etienne.fr/SEDiL>

Laurent Boyer<sup>(1)</sup>, Yann Esposito<sup>(1)</sup>, Amaury Habrard<sup>(2)</sup>, Jose Oncina<sup>(3)</sup> and Marc Sebban<sup>(1)</sup>



(1) Laboratoire Hubert Curien, UMR CNRS 5516, Université de Saint-Étienne, France (2) Laboratoire d'Informatique Fondamentale, UMR CNRS 6166, Aix-Marseille Université, France

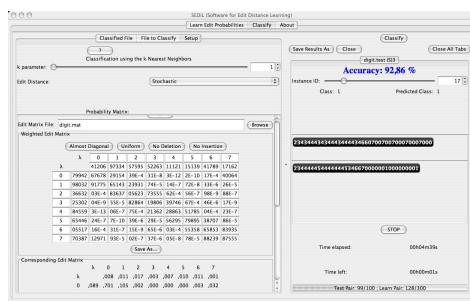
(3) Departamento de Lenguajes y Sistemas Informáticos, Universidad de Alicante, Spain



## Context

- \* Structured data (Strings, Trees, Graphs)
- \* Recurrent need of similarity measures in Machine Learning
  - ↳ Edit Distance (ED)
- \* Theoretical framework
  - ↳ Probabilistic models (Pair-HMMs, transducers, PDFAs, ...)
- \* Aim
  - ↳ To provide a platform grouping together the state of the art algorithms for learning edit parameters

⇒ SEDiL



Standard Edit Distance

## Definition

An **Edit script** is a sequence of edit operations (insertion, deletion, substitution) allowing the transformation of an input data X into an output data Y.

Probabilistic Edit Similarity

## Definition

The **Standard Edit Distance (ED)** between X and Y is the cost of the less costly edit script.

## Crucial step

Manual setting of the edit costs (depends on the domain).

## Example

```
pair of strings
input string a b
output string a a
```

Edit scripts = { (a → a, b → a); (a → λ, b → a), λ → a); (a → a, b → λ, λ → a); ... } where λ is the empty symbol

c	λ	a	b
λ	-	2	3
a	3	0	1
b	1	2	0

A priori fixed edit costs:

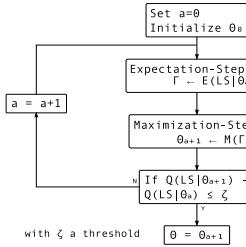
$$\text{ED}(a b, a a) = c(a \rightarrow a) + c(b \rightarrow a) = 0 + 2 = 2$$

$$\begin{aligned} p(a b, a a) &= \delta(a \rightarrow \lambda) \times \delta(b \rightarrow a) \times \delta(\lambda \rightarrow a) \\ &= 0.05 \times 0.2 \times 0.2 = 0.002 \\ \text{ES}_v(a b, a a) &= -\log(0.002) \\ &\approx 2.69 \end{aligned}$$

δ	λ	a	b
—	—	0.2	0.1
0.05	0.0	0.0	0.1
0.05	0.2	0.3	0.0

## ES learning algorithms in SEDiL

### Expectation-Maximization-based algorithms



### References

- [1] Ristani S., Yianilos P.: Learning string-edit distance. IEEE Transactions on Pattern Analysis and Machine Intelligence, 20(5). (1998), 522-532.
- [2] Durbin R., Eddy S., Krogh A., Mitchison G.: Biological sequence analysis. Cambridge University Press. (1998).
- [3] Oncina J., Sebban M.: Learning stochastic edit distance: application in handwritten character recognition. Pattern Recognition, 39(9). (2006), 1575-1587.
- [4] Bernard M., Habrard A., Sebban M.: Learning stochastic tree edit distance. ECML'06. (2006), 42-52.
- [5] Boyer L., Habrard A., Sebban M.: Learning metrics between tree structured data: Application to image recognition. ECML'07. (2007), 54-66.
- [6] Bernard M., Boyer L., Habrard A., Sebban M.: Learning probabilistic models of tree edit distance. Pattern Recognition, 41(8). (2008), 2611-2629.

### Features of the state of the art algorithms

[1]	[2]	[3]	[4,6]	[5]
strings	strings	strings	trees	trees
memoryless	non-memoryless	memoryless	memoryless	memoryless
Viterbi	stochastic	Viterbi	stochastic	stochastic
generative	generative	discriminative	generative	generative

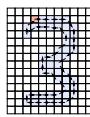
## Applications

Original data

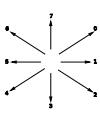
String representation

Tree representation

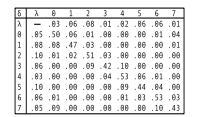
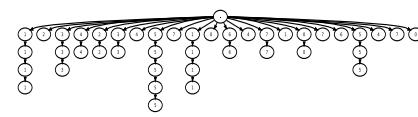
Learned parameters



Use of Freeman's codes



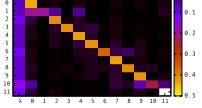
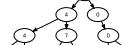
11112333442233455  
5557111066477100  
76555470



Handwritten Character Recognition



457200



Music Recognition