



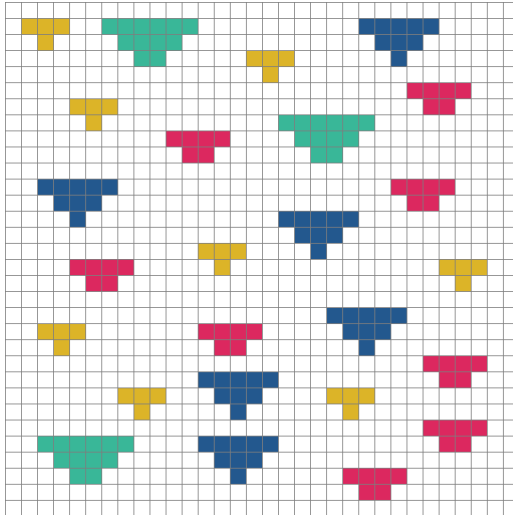
Geometric Pattern Mining using the MDL Principle

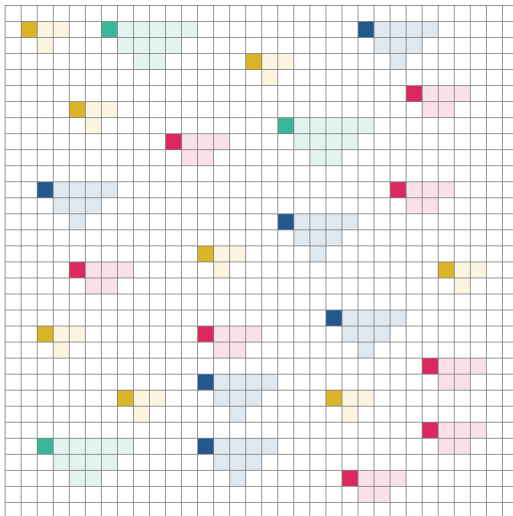
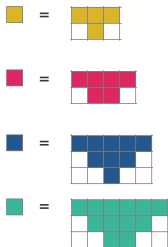
Micky Faas and Matthijs van Leeuwen

LIACS, Leiden University, Leiden, the Netherlands

April 20, 2020

[illegible]







Minimum Description Length

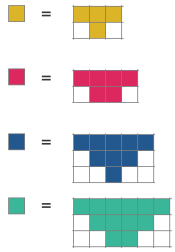
Idea: the best solution is a balance of model and instantiation complexity, given by:

$$\underbrace{L(H)}_{\text{Model}} + \underbrace{L(A|H)}_{\text{Data given model}}$$

In this case:

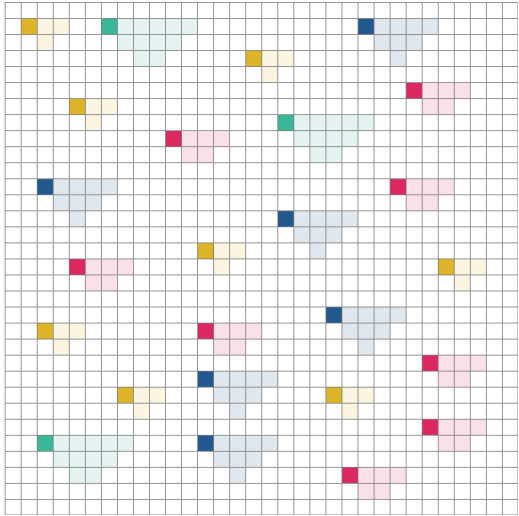
$$L \left(\underbrace{X = \begin{bmatrix} 1 & \cdot \\ \cdot & 1 \end{bmatrix}, Y = \begin{bmatrix} 1 \end{bmatrix}}_{\text{Model}} \right) + L \left(\underbrace{\begin{pmatrix} \text{Instantiation Matrix} \\ X & \cdot & \cdot & \cdot & \cdot & Y \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ X & \cdot & \cdot & \cdot & X & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ X & \cdot & X & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix}}_{\text{Data given model}} \right)$$

Model H

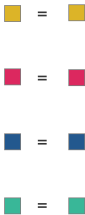


$L(H) = 887$
 $L(I) = 30$
 Compression: 89.5%

Instantiation Matrix I

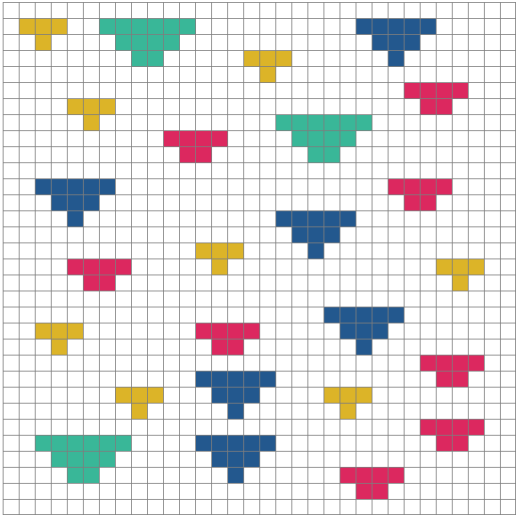


Model H

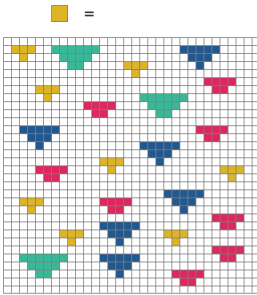


$L(H) = 4$
 $L(I) = 1024$
 Ratio: 100.4%

Instantiation Matrix I

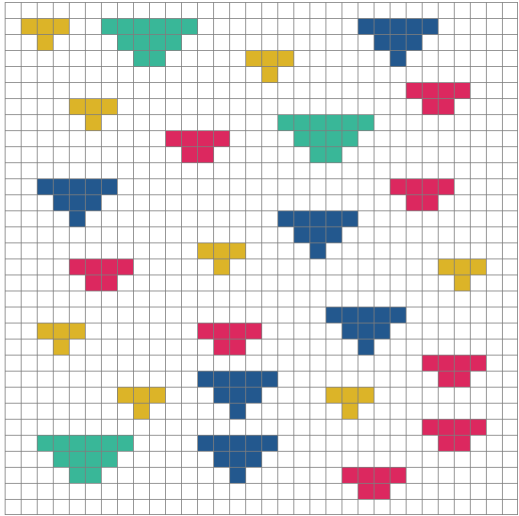


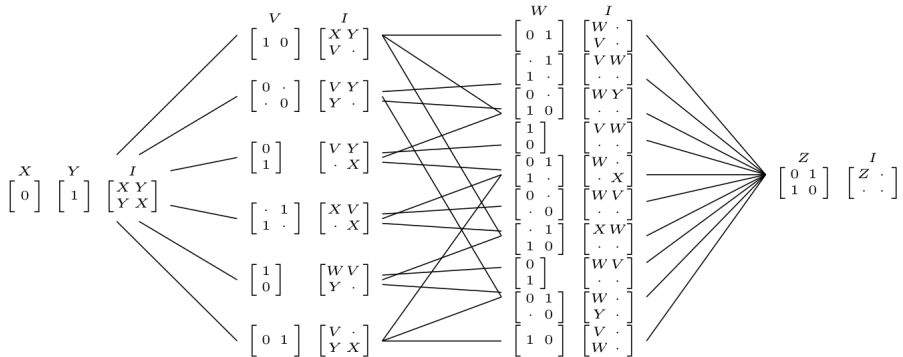
Model H



$L(H) = 1024$
 $L(I) = 1$
 Ratio: 100.1%

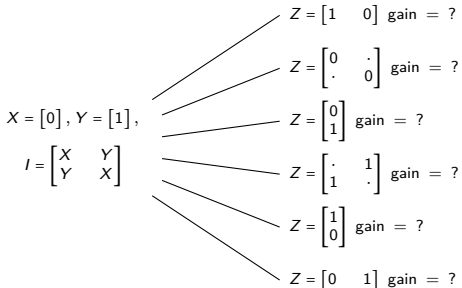
Instantiation Matrix I





Model space lattice for a 2×2 Boolean matrix. The V, W, and Z columns show which pattern is added in each step, while I depicts the current instantiation.

Candidates are generated by enumerating all combinations of two adjacent patterns that occur in the instantiation matrix.

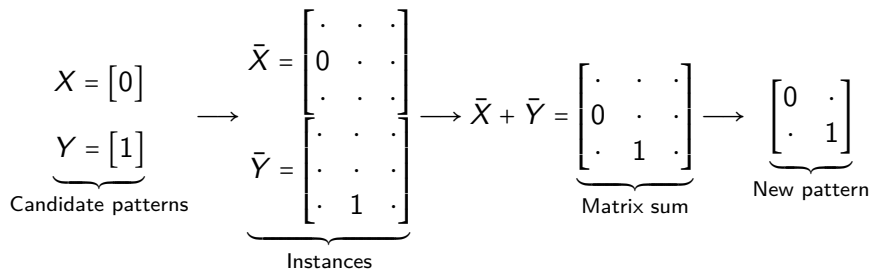


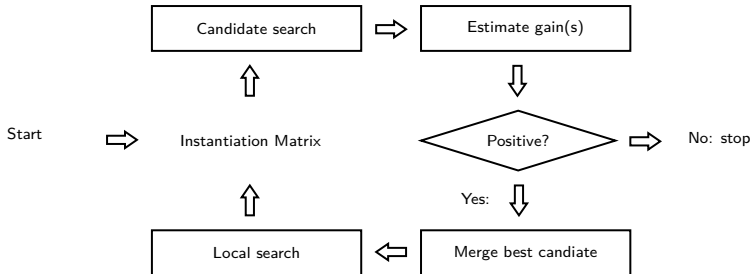
$$\underbrace{\Delta L(A', c)}_{\text{Gain}} = \underbrace{\left(L_1(H') + L_2(I') \right)}_{\text{New Lengths}} - \underbrace{\left(L_1(H) + L_2(I) \right)}_{\text{Old Lengths}}$$

L_1 and L_2 are independent length functions that compute the length of the model and the instantiation, respectively.

Please see the paper for more information on the encoding scheme.

We can construct complex patterns by repeatedly combining simpler ones. We use instances for this as they encode the position of one pattern relative to another.



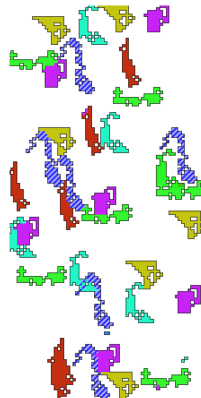




(a) Generated matrix



(b) Ground truth

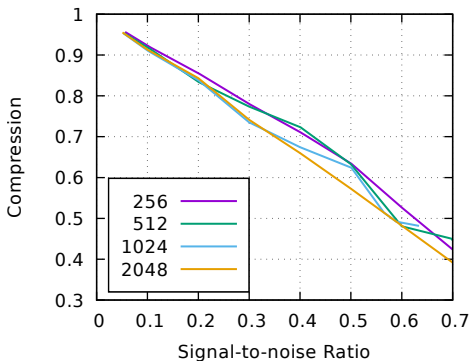


(c) Found patterns



(d) Difference

Figure: Synthetic patterns are added to a matrix filled with noise. The difference between the ground truth and the matrix reconstructed by the algorithm is used to compute precision and recall.



Results for different square matrices (256×256 to 1024×1024).
Signal-to-noise ratio is computed as $\frac{\text{signal}}{\text{signal} + \text{noise}}$.

The article was published at the Symposium on Intelligent Data Analysis (IDA) 2020

Archive link: <http://arxiv.org/abs/1911.09587>

Code repository: <https://github.com/mickymuis/libvouw>

Contact the authors:

Micky Faas <micky@educitty.org>

Matthijs van Leeuwen <m.van.leeuwen@liacs.leidenuniv.nl>

Thank you for watching!