

Geometric Formal Concept Analysis

- ① Abstract data cube category and lattice. Meet and join.
- ② Binding preliminaries. Maximal clusters.
- ③ Concrete data cube and description diagrams. Feature matrix, signatures, formal concepts.
- ④ The Formal Concept Analysis approach. Adjoint functors, Galois connection, closure, concept lattice.
- ⑤ Multiple concept binding. Tetrahedral and cubical points of view. Basic concepts.
- ⑥ Evaluation and selection of optimal description diagrams. Iterative binding algorithms.
- ⑦ Visualization of the final description. Decision tree row/column orders. Concept quality.
- ⑧ Further analysis of a description. Syntactic and semantic transfer.

Geometric Formal Concept Analysis

① Abstract data cube category and lattice

$C = [0,1]^n$ data cube

p_1, \dots, p_n coordinate functions on C ; properties

$$C_n := \{C\}$$

$$C_{n-1} := \{Z(p_i - v) \mid i=1, \dots, n, v=0,1\} \text{ generator facets; faces}$$

$$C_k := \left\{ \bigcap_{f \in G} f \mid G \subset C_{n-1}, |G|=n-k \right\} \setminus \{\emptyset\} \text{ k-dimensional facets}$$

$$\mathcal{C} := \left(\bigcup_{k=0}^{n-1} C_k, \text{ set inclusions} \right) \text{ facet category}$$

$$\psi: \mathcal{P}(C_{n-1}) \rightarrow \mathcal{P}(C) \quad G \mapsto \bigcap G \text{ geometric realization of symbolic facets}$$

$$A = \{G \in \mathcal{P}(C_{n-1}) \mid \psi G \neq \emptyset\}$$

$\psi|_A$ is injective.

$$\text{Obj } \mathcal{C} = \text{Image } \psi \setminus \{\emptyset\}$$

$$V := (\psi|_A)_* (\cup) \text{ meet in } \mathcal{C}$$

$$\wedge := (\psi|_A)_* (\cap) \text{ join in } \mathcal{C}$$

\mathcal{C} is a poset category.

In \mathcal{C} limits and colimits depend only on diagram objects; are products and coproducts.

$$F \subset \text{Obj } \mathcal{C}$$

$$\text{Proposition} \quad \lim F = VF = \sup F \quad \text{colim } F = \wedge F = \inf F$$

$$\left. \begin{array}{l} \text{Obj } \mathcal{C} \rightarrow C \quad f \mapsto \text{midpoint } f \\ \text{Mor } \mathcal{C} \rightarrow \text{Line segments in } C \end{array} \right\} \text{geometric realization of } \mathcal{C} \text{ as lattice}$$

② Binding preliminaries (compare [E.V., MES] "cluster" and "complex link")

$\tilde{\mathcal{C}}$ arbitrary category

$\text{Diag } \tilde{\mathcal{C}}$ diagrams in $\tilde{\mathcal{C}}$

$D_1, D_2 \in \text{Diag } \tilde{\mathcal{C}}$

$\{m_{ab}: a \rightarrow b\}_{a \in D_1, b \in D_2}$ a system of morphisms

$\{m_{ab}\}$ called a cluster if commute with D_1, D_2

cluster $\{m_{ab}\}$ called simple if factor through some $a' \rightarrow b'$.

Proposition Assume $\tilde{\mathcal{C}}$ is a poset category. Let $E, F \in \text{Obj } \tilde{\mathcal{C}}$

There is a cluster $E \rightarrow F$ if and only if $\text{colim } E \rightarrow \text{lim } F$.

It is automatically simple, and unique.

cluster $E \rightarrow F$ called maximal if $\text{colim } E \xrightarrow{\cong} \text{lim } F$

③ Concrete data cube and description diagrams

M $N \times n$ binary matrix table; observations or samples; feature matrix

S set of N rows of M ; signatures (assume non-redundant)

$S \subseteq C_0$

Assume $\Lambda S = C$. Otherwise reduce from C to ΛS by removing redundant properties.

Draw diagram in \mathcal{C} , one morphism $s \rightarrow (p, v)$ ($s \in S$, $(p, v) \in C_{n-1}$) for each entry of M .

D diagram in \mathcal{C}

D called a description (of S) if the subcategory of \mathcal{C} generated by D contains Draw

A (formal) concept, in the context of M or S , is a maximal cluster $E \rightarrow F$, $E \subseteq S$ and $F \subseteq C_{n-1}$.

Descriptions can be generated by binding concepts.

④ The Formal Concept Analysis approach

$$L := \psi^{-1} \circ \Lambda : \mathcal{P}(S) \rightarrow \mathcal{P}(C_{n-1})$$

$$R := V(-) \cap S : \mathcal{P}(C_{n-1}) \rightarrow \mathcal{P}(S)$$

L, R are adjoint functors, yielding Galois connection and closure operators:

$$E \mapsto \overline{E} = R(L(E)) \quad E \subseteq S$$

$$F \mapsto \overline{F} = L(R(F)) \quad F \subseteq C_{n-1}$$

theorem (Rudolph Wille) The concepts form a complete lattice. Projection to domain E or codomain F yields an isomorphism to lattice of closed subsets of S and closed subsets of C_{n-1} respectively. (The latter being contravariant).

theorem/algorithm () The closed sets (and hence the concepts) can be enumerated

by recursive closure of pairwise unions, starting from singletons in S .

theorem/algorithm ()

The closed sets can be enumerated in "lectic order" with respect to ordering of $\{p_1, \dots, p_n\}$

The FCA approach is concerned with just one description of S : the concept lattice.

⑤ Multiple concept binding; tetrahedral and cubical points of view

$$e := |E| \quad f := |F|$$

Binding $E \rightarrow F$ changes #edges: $ef \rightarrow e+f$, #nodes: $+1$ (when starting from Draw)

If two concepts are completely independent, their edge/node savings are also independent (they add).

$E \rightarrow F, E' \rightarrow F'$ two concepts (say, c and c')

$$e_1 := |E \setminus E'| \quad e_2 := |E \cap E'| \quad e_3 := |E' \setminus E|$$

$$f_1 := |F \setminus F'| \quad f_2 := |F \cap F'| \quad f_3 := |F' \setminus F|$$

$$e := e_1 + e_2 + e_3 = |E \cup E'|$$

$$f := f_1 + f_2 + f_3 = |F \cup F'|$$

$u := E \cup E' \rightarrow F \cap F'$ upper concept induced by c, c'

$l := E \cap E' \rightarrow F \cup F'$ lower concept induced by c, c'

Binding c, c', u, l changes #edges: $\underbrace{\sum_{i=1}^3 e_i f_i + \sum_{j=1}^2 (e_j f_{j+1} + e_{j+1} f_j)}_{ef - (e_1 f_3 + e_3 f_1)} \rightarrow e + f + 4$, #nodes: $+4$

General case, K concepts c_1, \dots, c_K

If all intersections are non-empty, the concepts induced by $\{c_i\}$ are labelled by the vertices of the K -dimensional cube; equivalently, by the facets of the $(K-1)$ -dimensional simplex.

Aside from morphisms emanating directly from S or C_{n-1} , the morphisms of the bound diagram are labelled by the 1-skeleton of the K -cube.

If some intersections are empty, one needs the quotient graph by the edges emanating from the nodes labelling empty intersection.

Problem Identify appropriately minimal collection of concepts which induce the rest.

Attempted solution.

$E \rightarrow F$ called basic if one of the following equivalent conditions hold:

- E is neither the union of two proper closed subsets nor the intersection of two proper closed supersets.

- F is

- E is not the union of two proper closed subsets and F is not the union of two proper closed subsets.

- E is not the intersection of two proper closed supersets and F is not the intersection of two proper closed supersets.

Problem Identify all basic concepts with an algorithm.

⑥ Evaluation and selection of optimal description diagrams

Complexity metrics/functions on set of descriptions

- total number of edges
- total number of edges + total number of nodes
- total edge length in 2D embedding (Steiner problem)

~~Algorithmic~~ strategies

Algorithm initialization strategies

- Sequence. Start from a deterministic pre-defined sequence of concepts to be bound, e.g. lexic order.
- Basis. Start from naive simultaneous binding of a pre-defined finite set of core concepts, e.g. basic concepts.

Bind until incremental improvement of complexity metric falls below error threshold, or until next step is not an improvement.

In basis case, should consider intersection/union of all previously bound concepts, or else a level-based method:

Level 1. Bind all pairs intersection/unions of basis until improvement threshold, replacing the pair with the new, intersection and union. (Still K concepts on the dynamic list).

Level 2. Bind all triples intersection/unions of basis until improvement threshold, replacing the 3 with the 2 new, intersection and union.

... termination is likely since the number of concepts on deck after Level 1 is decreasing with each increment.

⑦ Visualization of the final description

Append concept membership vectors to M as new features, use decision tree w.r.t. them to order columns of M .

Append concept membership vectors to M as quasi-samples, use decision tree w.r.t. them to order rows of M .

Use opacity map to display "quality" of description per entry of M , e.g.

- the size (area) of the largest concept capturing that entry
- average size (area) of the concepts capturing that entry
- as above, assessing concept quality differently, e.g. variance of feature values pre-dichotomization, in case M was created by dichotomization from continuous values.

⑧ Further analysis of a description: syntax/semantics and transfer

D a description of S

α a subdiagram of D containing C_{n-1} and none of S (α for "abstract")

K a subdiagram of D containing S and none of C_{n-1} (K for "concrete")

A splitting of D is α and K such that D is the union of α , K , and any morphisms sourced in K and targeted in α which belong to D .

Transfer to a new formal context M'/S' with at least the same feature set as M/S ,

means the description D' of S' obtained by binding all closures of the upper sets of each object of α .

Transfer may be used within a single dataset to set up cross-validation.

The above may be called syntactic transfer.

A transfer procedure involving K may be called semantic transfer:

Extend to a new formal context M'/S' with more features on the same sample set S , e.g. outcome features previously withheld. Calculate D' by closure of each lower set of an object of K in the new context.