

Categorical Foundation of Explainable AI

J. Eickhoff

Lecture: Explainable Artificial Intelligence, 10/20/2025

Inhalt

- 1 Category Theory
- 2 Institutional model theory
- 3 Special Categories
- 4 Categorical Framework of Explainable AI
- 5 Impact on XAI

"Mathematics is the art of giving the same name to different things."

– Henri Poincaré, Science et Méthode (1908)

Category Theory

Category Theory provides a unifying mathematical language for describing structures and relationships between them.

It abstracts the essential idea of *composition* and *identity* found across all areas of mathematics.

Idea

"Mathematics is about structures, and category theory is the study of structure itself."

A *category* consists of **objects** and **morphisms** (arrows) between them, with composition and identity satisfying associativity and unity laws.

Category Theory

Definition (Category)

A *category* \mathcal{C} consists of the following components:

- a class of **objects**, denoted by $\text{Ob}(\mathcal{C})$,
- for each pair of objects $X, Y \in \text{Ob}(\mathcal{C})$, a class of **morphisms**

$$\text{Hom}_{\mathcal{C}}(X, Y),$$

which includes a *identity morphism* $\text{Id}_X \in \text{Hom}_{\mathcal{C}}(X, X)$ for each object X ,

- and a **composition** operation

$$\circ: \text{Hom}_{\mathcal{C}}(X, Y) \times \text{Hom}_{\mathcal{C}}(Y, Z) \rightarrow \text{Hom}_{\mathcal{C}}(X, Z),$$

satisfying

$$\text{Id}_Y \circ f = f = f \circ \text{Id}_X, \quad (h \circ g) \circ f = h \circ (g \circ f).$$

Definition (Opposite Category)

Let \mathcal{C} be a category. Then \mathcal{C}^{op} is the *opposite category* of \mathcal{C} with following data:

- $\text{Ob}(\mathcal{C}^{\text{op}}) = \text{Ob}(\mathcal{C})$
- $\text{Hom}_{\mathcal{C}^{\text{op}}}(X, Y) = \text{Hom}_{\mathcal{C}}(Y, X)$

Remark: We call a category *small* if the objects and morphisms form sets. In case the class of morphisms form sets, we call the category *locally small*.

Definition (Functor)

A functor $F: \mathcal{C} \rightarrow \mathcal{D}$ assigns:

- to each object $X \in \mathcal{C}$ an object $F(X) \in \mathcal{D}$,
- to each morphism $f: X \rightarrow Y$ in \mathcal{C} a morphism

$$F(f): F(X) \rightarrow F(Y)$$

in \mathcal{D} ,

such that:

- $F(\text{Id}_X) = \text{Id}_{F(X)}$ for every object X ,
- $F(g \circ f) = F(g) \circ F(f)$ for all composable morphisms f, g .

Definition (Natural Transformation)

Let $F, G : \mathcal{C} \rightarrow \mathcal{D}$ be functors. A *natural transformation* $\eta : F \Rightarrow G$ assigns to each object $X \in \mathcal{C}$ a morphism

$$\eta_X : F(X) \rightarrow G(X)$$

in \mathcal{D} such that for every morphism $f : X \rightarrow Y$ in \mathcal{C} the following diagram commutes:

$$\begin{array}{ccc} F(X) & \xrightarrow{F(f)} & F(Y) \\ \eta_X \downarrow & & \downarrow \eta_Y \\ G(X) & \xrightarrow{G(f)} & G(Y) \end{array}$$

If each η_X is an isomorphism, then η is called a *natural isomorphism*.

Definition (Equivalence of Categories)

Two categories \mathcal{C} and \mathcal{D} are said to be *equivalent* if there exist functors

$$F : \mathcal{C} \rightarrow \mathcal{D} \quad \text{and} \quad G : \mathcal{D} \rightarrow \mathcal{C}$$

such that the compositions

$$G \circ F \cong \text{Id}_{\mathcal{C}} \quad \text{and} \quad F \circ G \cong \text{Id}_{\mathcal{D}}$$

are naturally isomorphic to the respective identity functors.

Inhalt

- 1 Category Theory
- 2 Institutional model theory**
- 3 Special Categories
- 4 Categorical Framework of Explainable AI
- 5 Impact on XAI

Institutional model theory

Institutional Model Theory is an abstract framework for logic and model theory.

It provides a *unified categorical description* of various logical systems (e.g. first-order logic, modal logic, algebraic specification).

Idea

"Truth is invariant under change of notation."

That is, the relation between syntax (sentences) and semantics (models) remains consistent across different signatures.

Signature

Definition (Signature)

A *signature* is a quadruple

$$\Sigma = (\mathcal{S}_{\text{func}}, \mathcal{S}_{\text{rel}}, \mathcal{S}_{\text{const}}, \text{ar}),$$

where

- $\mathcal{S}_{\text{func}}$ and \mathcal{S}_{rel} are disjoint sets of
 - ▶ function symbols (e.g., $+$, \times),
 - ▶ relation symbols or predicates (e.g., \leq , \in),
 - ▶ constant symbols (e.g., 0 , 1),
- and $\text{ar} : \mathcal{S}_{\text{func}} \cup \mathcal{S}_{\text{rel}} \rightarrow \mathbb{N}$ assigns a natural number (called the *arity*) to each function or relation symbol.

A function or relation symbol is called *n-ary* if its arity is n .

Institutional Model Theory

Let Cat^{op} denote the opposite category of small categories. Consider

- a category Sign of *signatures*;
- a functor $\text{Sen} : \text{Sign} \rightarrow \text{Set}$ assigning to each Σ its sentences $\text{Sen}(\Sigma)$;
- a functor $\text{Mod} : \text{Sign} \rightarrow \text{Cat}^{\text{op}}$ assigning to each Σ its category of models $\text{Mod}(\Sigma)$;
- a satisfaction relation

$$\models_{\Sigma} \subseteq |\text{Mod}(\Sigma)| \times \text{Sen}(\Sigma).$$

Institution Model Theory

Definition (Institution)

An **institution** is a quadruple

$$(\text{Sign}, \text{Sen}, \text{Mod}, \models).$$

For a signature morphism $\sigma : \Sigma \rightarrow \Sigma'$:

- **sentence translation** $\text{Sen}(\sigma) : \text{Sen}(\Sigma) \rightarrow \text{Sen}(\Sigma')$,
- **reduct functor** $\text{Mod}(\sigma) : \text{Mod}(\Sigma') \rightarrow \text{Mod}(\Sigma)$ (contravariant via Cat^{op}).

Satisfaction condition: for all $M' \in \text{Mod}(\Sigma')$ and $\varphi \in \text{Sen}(\Sigma)$,

$$M' \models_{\Sigma'} \text{Sen}(\sigma)(\varphi) \iff \text{Mod}(\sigma)(M') \models_{\Sigma} \varphi.$$

Example: Institution of Propositional Logic

Propositional Logic can be represented as an institution where:

- **Signatures:** $\Sigma = \{a, b, c, \neg, \wedge, \vee, \rightarrow, \top, \perp\}$ with arity:
 $\text{ar}(\neg) = 1, \quad \text{ar}(\wedge) = \text{ar}(\vee) = \text{ar}(\rightarrow) = 2, \quad \text{ar}(a) = \text{ar}(b) = \text{ar}(c) = 0$
- **Sentences:** $\text{Sen}(\Sigma)$ – all Boolean formulas over Σ
- **Models:** $\text{Mod}(\Sigma) = \{0, 1\}^{\{a, b, c\}}$
- $M \models_{\Sigma} \varphi$ iff φ evaluates to true under M

Idea

Truth depends only on Boolean combinations of fixed atomic propositions.

Explicit Example

For $M(a) = 1, M(b) = 1, M(c) = 0$ we have:

$$(a \wedge b) \rightarrow \neg c = (1 \wedge 1) \rightarrow \neg 0 = 1 \rightarrow 1 = 1.$$

Hence $M \models_{\Sigma} \varphi$.

Example: Institution of First-Order Logic

First-Order Logic extends propositional logic with variables, quantifiers, and structure.

- Sign – first-order signatures with functions and relations (e.g. $+$, 0 , $<$),
- $\text{Sen}(\Sigma)$ – well-formed formulas over Σ (e.g. $\forall x (x + 0 = x)$),
- $\text{Mod}(\Sigma)$ – Σ -structures interpreting symbols,
e.g. $(\mathbb{N}, +, 0, <)$,
- $M \models_{\Sigma} \varphi$ iff φ is true in M under all variable assignments.

Idea

Truth now depends on the interpretation of symbols in a structure

Explicit Example

In $M = (\mathbb{N}, +, 0, <)$:

$$\forall x (x + 0 = x)$$

means “for all natural numbers x , $x + 0 = x$ ”. This holds for every $x \in \mathbb{N}$,
so $M \models_{\Sigma} \varphi$.

Using Structure in First-Order Logic

In first-order logic, models have internal structure:

$$M = (|M|, f^M, R^M, c^M, \dots)$$

This allows statements that talk about elements and relations inside M .

In the institutional view:

- $M \in \text{Mod}(\Sigma)$ abstracts away from internal details.
- Only the satisfaction relation $M \models_{\Sigma} \varphi$ is required.

This abstraction allows a unified treatment of different logics.

Using Structure in First-Order Logic

Example

Let $\Sigma = \{+, 0, <\}$ and $M = (\mathbb{N}, +, 0, <)$.

We can express:

- **Addition property:**

$$\forall x \forall y \forall z (x < y \rightarrow x + z < y + z)$$

- **Identity law:**

$$\forall x (x + 0 = x)$$

- **Existence of successor:**

$$\forall x \exists y (x < y)$$

Each statement uses the *internal structure* of $(\mathbb{N}, +, 0, <)$: functions (+), constants (0), and relations (<).

Example: Institution of Equational Logic

Equational Logic provides an algebraic institution based on equations between terms.

- Sign – operations and sorts, e.g. $\{+, 0, -\}$ for groups,
- $\text{Sen}(\Sigma)$ – all equations $t_1 = t_2$ built from Σ -terms,
- $\text{Mod}(\Sigma)$ – Σ -algebras interpreting the operations,
- $A \models_{\Sigma} (t_1 = t_2)$ iff both terms evaluate equally in A .

Idea

Equational logic generalizes algebra: models are algebras, and truth means equality.

Explicit Example

In $A = (\mathbb{Z}, +, -, 0)$:

$$x + 0 = x$$

is true for all $x \in \mathbb{Z}$ since $x + 0 = x$ in integer addition.

Hence $A \models_{\Sigma} (x + 0 = x)$.

Example: The Real Numbers

The reals \mathbb{R} can be viewed both algebraically and logically.

Let the signature be

$$\Sigma = \{+, 0, <\}.$$

A model for this signature can be written as

$$M = (\mathbb{R}, +^M, <^M, 0^M, 1^M, \cdot^M),$$

where:

- $|M| = \mathbb{R}$ is the carrier set,
- $+^M$ is the usual addition on \mathbb{R} ,
- $<^M$ is the usual order relation on \mathbb{R} ,
- $0^M = 0$ is a constant,
- $1^M = 1$ is constant,
- \cdot^M is the usual multiplication on \mathbb{R} .

Example: The Real Numbers

- **As an algebra (Equational Logic):**

$$A = (\mathbb{R}, +, \cdot, 0, 1)$$

with equations such as

$$x + 0 = x, \quad x \cdot 1 = x, \quad x + y = y + x.$$

Truth means both sides evaluate equally in A .

- **As a structure (First-Order Logic):**

$$M = (\mathbb{R}, +, \cdot, 0, 1, <)$$

allowing statements like

$$\forall x \forall y (x < y \rightarrow x + 1 < y + 1).$$

Idea

Adding the relation $<$ turns the algebra into a first-order structure.

Example: Changing Between Models

In an institution, a signature morphism

$$\sigma : \Sigma \rightarrow \Sigma'$$

induces a functor

$$\text{Mod}(\sigma) : \text{Mod}(\Sigma') \rightarrow \text{Mod}(\Sigma)$$

that *reduces* a model by forgetting structure.

Example (from richer to simpler model)

Let

$$\Sigma = \{+, 0\}, \quad \Sigma' = \{+, 0, <\},$$

and σ forgets the relation $<$.

Then:

$$\text{Mod}(\Sigma') \ni M' = (\mathbb{R}, +, 0, <) \mapsto \text{Mod}(\sigma)(M') = (\mathbb{R}, +, 0) \in \text{Mod}(\Sigma).$$

Hence the first-order structure $(\mathbb{R}, +, 0, <)$ is *reduced* to its algebraic part $(\mathbb{R}, +, 0)$.

Example: Translating Sentences Between Signatures

In an institution, a signature morphism

$$\sigma : \Sigma \rightarrow \Sigma'$$

induces:

$$\text{Sen}(\sigma) : \text{Sen}(\Sigma) \rightarrow \text{Sen}(\Sigma'),$$

which *translates sentences* by renaming or extending symbols.

Example: Translating Sentences Between Signatures

Example (formula translation)

Let

$$\Sigma = \{+, 0\}, \quad \Sigma' = \{+, 0, <\},$$

and define σ as the inclusion map.

Then for a formula

$$\varphi = \forall x (x + 0 = x) \in \text{Sen}(\Sigma)$$

we get

$$\text{Sen}(\sigma)(\varphi) = \forall x (x + 0 = x) \in \text{Sen}(\Sigma'),$$

the same statement, but now expressed in the richer language that also allows the relation $<$.

Idea

$\text{Sen}(\sigma)$ translates or embeds formulas so that truth remains preserved when the language changes.

Inhalt

- 1 Category Theory
- 2 Institutional model theory
- 3 Special Categories**
- 4 Categorical Framework of Explainable AI
- 5 Impact on XAI

Monoidal Category

Definition

A *monoidal category* is a category \mathcal{C} equipped with

- a bifunctor $\otimes : \mathcal{C} \times \mathcal{C} \rightarrow \mathcal{C}$,
- a unit object $I \in \mathcal{C}$,
- natural isomorphisms

$$\alpha_{A,B,C} : (A \otimes B) \otimes C \xrightarrow{\cong} A \otimes (B \otimes C)$$

$$\lambda_A : I \otimes A \xrightarrow{\cong} A, \quad \rho_A : A \otimes I \xrightarrow{\cong} A$$

for all $A, B, C \in \mathcal{C}$, satisfying the *pentagon* and *triangle* coherence axioms (Mac Lane's coherence theorem).

Monoidal Category

Pentagon Diagramm

$$\begin{array}{ccccc} ((A \otimes B) \otimes C) \otimes D & \xrightarrow{\alpha_{A,B,C} \otimes \text{id}_D} & (A \otimes (B \otimes C)) \otimes D & \xrightarrow{\alpha_{A,B \otimes C,D}} & A \otimes ((B \otimes C) \otimes D) \\ \downarrow \alpha_{A \otimes B,C,D} & & & & \downarrow \text{id}_A \otimes \alpha_{B,C,D} \\ (A \otimes B) \otimes (C \otimes D) & \xrightarrow{\alpha_{A,B,C \otimes D}} & & & A \otimes (B \otimes (C \otimes D)) \end{array}$$

Triangle Diagramm

$$\begin{array}{ccc} (A \otimes I) \otimes B & \xrightarrow{\alpha_{A,I,B}} & A \otimes (I \otimes B) \\ \searrow \rho_A \otimes \text{id}_B & & \swarrow \text{id}_A \otimes \lambda_B \\ & A \otimes B & \end{array}$$

Symmetric Monoidal Category

Definition (Symmetric Monoidal Category)

A *symmetric monoidal category* is a monoidal category $(\mathcal{C}, \otimes, I, \alpha, \lambda, \rho)$ together with a natural isomorphism (the *symmetry*)

$$s_{A,B} : A \otimes B \xrightarrow{\cong} B \otimes A$$

such that the following diagrams commute for all $A, B, C \in \mathcal{C}$:

Symmetric Monoidal Category

Unit coherence

$$\begin{array}{ccc} A \otimes I & \xrightarrow{s_{A,I}} & I \otimes A \\ & \searrow \rho_A \quad \swarrow \lambda_A & \\ & A & \end{array}$$

Associativity coherence

$$\begin{array}{ccccc} (A \otimes B) \otimes C & \xrightarrow{s_{A,B} \otimes \text{id}_C} & (B \otimes A) \otimes C & \xrightarrow{\alpha_{B,A,C}} & B \otimes (A \otimes C) \\ \alpha_{A,B,C} \downarrow & & & & \downarrow \text{id}_B \otimes s_{A,C} \\ A \otimes (B \otimes C) & \xrightarrow{s_{A,B \otimes C}} & (B \otimes C) \otimes A & \xrightarrow{\alpha_{B,C,A}} & B \otimes (C \otimes A) \end{array}$$

Inverse law

$$\begin{array}{ccc} A \otimes B & \xrightarrow{s_{A,B}} & B \otimes A \\ & \searrow \text{id}_{A \otimes B} \quad \swarrow s_{B,A} & \\ & A \otimes B & \end{array}$$

Example: The Category \mathbf{Vect}_k

\mathbf{Vect}_k is the category of vector spaces over a field k :

- **Objects:** vector spaces V, W, \dots
- **Morphisms:** k -linear maps $f : V \rightarrow W$

Monoidal Structure

- **Tensor product:** $V \otimes W$
- **Unit object:** k (the base field)
- **Associator:**

$$\alpha_{U,V,W} : (U \otimes V) \otimes W \rightarrow U \otimes (V \otimes W)$$

- **Unit isomorphisms:** $\lambda_V : k \otimes V \rightarrow V, \quad \rho_V : V \otimes k \rightarrow V$

Example: The Category \mathbf{Vect}_k

Symmetry

$$s_{V,W} : V \otimes W \longrightarrow W \otimes V, \quad v \otimes w \longmapsto w \otimes v$$

- **Involution:** $s_{W,V} \circ s_{V,W} = \text{id}_{V \otimes W}$
- **Naturality:** $(g \otimes f) \circ s_{V,W} = s_{V',W'} \circ (f \otimes g)$

$(\mathbf{Vect}_k, \otimes, k, \alpha, \lambda, \rho, s)$ is a **symmetric monoidal category**.

Tensoring combines vector spaces, and the swap map shows that the order of tensor factors does not matter up to canonical isomorphism.

Definition (Cartesian Monoidal Category)

Definition (Cartesian Monoidal Category)

A **cartesian monoidal category** is a monoidal category

$$(\mathcal{C}, \times, 1, \alpha, \lambda, \rho)$$

whose monoidal product \times is given by the **categorical product**, and whose unit object is the **terminal object** 1.

Remark:

- The **categorical product** $A \times B$ has projections $\pi_1 : A \times B \rightarrow A$, $\pi_2 : A \times B \rightarrow B$, and any pair (f, g) factors uniquely through $A \times B$.
- The **terminal object** 1 is an object with exactly *one morphism* $A \rightarrow 1$ for every A .

Example: The Category $(\text{Set}, \times, 1)$

Structure

- **Objects:** Sets A, B
- **Morphisms:** Functions $f : A \rightarrow B$
- **Monoidal product:** Cartesian product $A \times B$
- **Unit:** Singleton set $1 = \{*\}$ (terminal object)

Canonical maps

- Projections: $\pi_1(a, b) = a, \pi_2(a, b) = b$
- Pairing: $\langle f, g \rangle(x) = (f(x), g(x))$

The product $A \times B$ satisfies the universal property, and every set A has a unique map $A \rightarrow 1$. Thus $(\text{Set}, \times, 1)$ is a **cartesian monoidal category**.

Co- vs. Cartesian Monoidal Categories

- $(\text{Set}, \times, 1) \rightarrow$ **cartesian monoidal** (product of sets, terminal object $1 = \{*\}$)
- $(\text{Vect}_k, \otimes, k) \rightarrow$ **monoidal, not cartesian** (no canonical projections $V \otimes W \rightarrow V, W$)
- $(\text{Vect}_k, \oplus, 0) \rightarrow$ **co-cartesian monoidal** (direct sum, trivial vector space as unit)

Feedback Monoidal Category

Definition (Feedback Monoidal Category)

A **feedback monoidal category** is a **symmetric monoidal category** equipped with:

- an **endofunctor** $F : \mathcal{C} \rightarrow \mathcal{C}$,
- and for all $X, Y, S \in \mathcal{C}$ an operation

$$\circlearrowleft_S : \text{Hom}(X \otimes F(S), Y \otimes S) \longrightarrow \text{Hom}(X, Y),$$

called the **feedback operator**.

This operator must satisfy the following **axioms** for all suitable f, g, h :

- **(Tightening)** $\circlearrowleft_S ((h \otimes 1_S) \circ f \circ (g \otimes 1_{FS})) = h \circ \circlearrowleft_S (f) \circ g$
- **(Joining)** $\circlearrowleft_{S \otimes T} (f) = \circlearrowleft_S (\circlearrowleft_T (f))$
- **(Vanishing)** $\circlearrowleft_I (f) = f$
- **(Strength)** $\circlearrowleft_S (f \otimes g) = g \otimes \circlearrowleft_S (f)$
- **(Sliding)** $\circlearrowleft_T ((1_Y \otimes g) \circ f) = \circlearrowleft_S (f \circ (1_X \otimes g))$

Feedback Monoidal Categories

Idea: Feedback monoidal categories describe **processes with feedback loops** in a monoidal setting.

- The operator \odot_S connects an output back to an input, forming a **feedback loop**.
- The axioms ensure feedback behaves coherently with composition (\circ) and tensor product (\otimes).

Remark: If the monoidal structure is **cartesian**, we call it a **feedback cartesian category**. Such categories model stateful systems or causal stream functions.

Cartesian Streams

Definition (Cartesiann Streams)

A **cartesian stream** $f : \chi \rightarrow \psi$ consists of two sequences of sets

$$\chi = (X_0, X_1, \dots), \quad \psi = (Y_0, Y_1, \dots)$$

and, for each $n \geq 0$, a map

$$f_n : X_n \times \dots \times X_0 \longrightarrow Y_n.$$

For each n , define the cumulative output

$$\hat{f}_n : X_n \times \dots \times X_0 \longrightarrow Y_n \times \dots \times Y_0$$

recursively by

$$\hat{f}_0 := f_0, \quad \hat{f}_{n+1} := (f_{n+1} \times \hat{f}_n) \circ (1_{X_{n+1}} \times \nu_{X_n \times \dots \times X_0}),$$

where $\nu_Z : Z \rightarrow Z \times Z$ is the diagonal (duplication) map.

Cartesian Streams

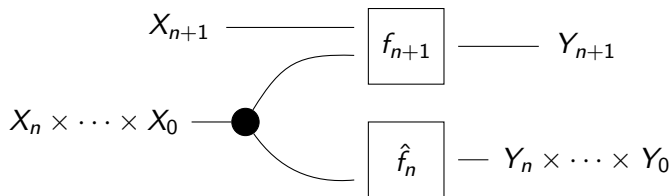
Composition of streams: for $f : \chi \rightarrow \psi$ and $g : \psi \rightarrow \zeta$,

$$(g \circ f)_n := g_n \circ \hat{f}_n.$$

Identity stream: $(1_\chi)_n := \pi_{\chi_n}$ (the n -th projection).

Connection: Cartesian streams form a category $\text{Stream}_{\text{Set}}$ that models *stateful or causal processes*. This category is a special case of a **feedback cartesian monoidal category**, where the cartesian product serves as the tensor, and the feedback operation corresponds to connecting the output of one step as input to the next.

String Diagrams



Free Category

Definition (Free Category)

Let $G = (V, E, s, t)$ be a directed graph with vertex set V , edge set E , and source and target maps $s, t : E \rightarrow V$.

The *free category* on G , denoted $F(G)$, is defined by:

- $\text{Ob}(F(G)) = V$;
- $\text{Hom}_{F(G)}(v, w)$ consists of all finite composable paths (e_n, \dots, e_1) with $s(e_1) = v$, $t(e_n) = w$, together with the empty path $()$;
- composition is concatenation of paths, and $\text{id}_v = ()$.

The assignment

$$F : \text{Digraph} \rightarrow \text{Cat}, \quad G \mapsto F(G),$$

is the *free category functor* and Cat is the category of small categories.

Inhalt

- 1 Category Theory
- 2 Institutional model theory
- 3 Special Categories
- 4 Categorical Framework of Explainable AI**
- 5 Impact on XAI

Categorical Framework of Explainable AI

- Learning is viewed as an **iterative process with feedback**.
- The process involves a **model/explainer** function that updates internal parameters using feedback from the environment (via an optimizer).
- Formally, an abstract learning agent is a **morphism** in the free feedback Cartesian monoidal category **XLearn**.

Categorical Framework of Explainable AI

XLearn is generated by:

- **Objects:** X, Y, Y^*, P, E
 - ▶ Represent input, output, supervision, parameter, and explanation types.
- **Model/Explainer morphism:**

$$\eta : X \times P \rightarrow Y \times E$$

- ▶ Produces predictions in Y and explanations in E .

- **Optimizer morphism:**

$$\nabla_Y : Y^* \times Y \times P \rightarrow P$$

- ▶ Updates parameters in P .
 - ▶ Uses supervisions in Y^* , model predictions in Y , and current parameters in P .

Definition (XLearn)

XLearn is the free feedback cartesian category generated by the objects

$$X, Y, Y^*, P, E$$

and by the morphisms

$$\eta : X \times P \rightarrow Y \times E \quad \text{and} \quad \nabla_Y : Y^* \times Y \times P \rightarrow P.$$

Abstract Learning Agent

Remark: In a feedback cartesian monoidal category, each object X is equipped with morphisms

$$\nu_X : X \rightarrow X \times X \quad \text{and} \quad \epsilon_X : X \rightarrow e$$

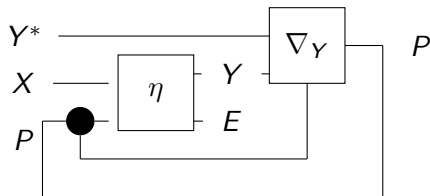
that make it possible to *copy* and *discard* objects.

Definition (Abstract Learning Agent)

An *abstract learning agent* is the morphism in $XLearn$ given by:

$$\circlearrowleft_P \left(\nabla_Y \circ (1_{Y^*} \times Y \times \epsilon_E) \circ (1_{Y^*} \times \eta \times 1_P) \circ (1_{Y^*} \times X \times \nu_P) \right).$$

Abstract Learning Agent



Concrete Learning and Explaining Agents

- The free category *XLearn* captures key features of learning agents at an abstract level.
- We can instantiate these agents in *concrete forms* using a **feedback functor** from *XLearn* to the category of Cartesian streams over **Set**, called *Stream_{Set}*.
- This functor maps the abstract structure to concrete settings involving:
 - ▶ Various explainers (e.g., decision trees, logistic regression),
 - ▶ Different input data types (e.g., images, text),
 - ▶ Supervisions, outputs, parameters, and explanations.
- Establishing this mapping requires defining a specific functor called the **translator**.

Agent Translator

Definition (Agent Translator)

An *agent translator* is a feedback cartesian functor

$$\mathcal{T} : \text{XLearn} \rightarrow \text{Stream}_{\text{Set}}.$$

Classes of Translators

- Among translators, we distinguish two main classes:
 - ▶ **Learning translators:** instantiate *learning agents*.
 - ▶ **Explaining translators:** instantiate *explainable learning agents*.
- Intuitively:
 - ▶ A **learning agent** is an instance of an abstract learning agent that does *not* provide explanations.
 - ▶ A **explaining learning agent** outputs a *non-empty explanation*.

Learning Agent

Definition (Learning Agent)

Given an agent translator \mathcal{T} with

$$\mathcal{T}(E) = \{*\}^{\mathbb{N}},$$

where $\{*\}$ is a singleton set.

A *learning agent* (*LA*) is the image $\mathcal{T}(\alpha)$, where α is the abstract learning agent.

Explaining Learning Agent

Definition (Explaining Learning Agent)

Let \mathcal{T} be an agent translator, I an institution, and α the abstract learning agent.

The image $\mathcal{T}(\alpha)$ is called

- a *syntactic explaining learning agent* if

$$\mathcal{T}(E) = \text{Sen}(\Sigma)^{\mathbb{N}},$$

- and a *semantic explaining learning agent* if

$$\mathcal{T}(E) = \text{Mod}(\Sigma)^{\mathbb{N}},$$

for some signature Σ of I .

Remark: Concrete instances of both syntactic and semantic explaining learning agents are referred to as *explaining learning agents (XLA)*.

Inhalt

- 1 Category Theory
- 2 Institutional model theory
- 3 Special Categories
- 4 Categorical Framework of Explainable AI
- 5 Impact on XAI**

Simplified Notation

To simplify the notation, we use the following shortcuts:

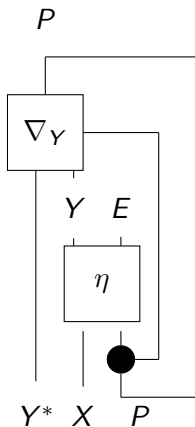
$$\mathcal{X} = \mathcal{T}(X), \quad \mathcal{Y} = \mathcal{T}(Y), \quad \mathcal{Y}^* = \mathcal{T}(Y^*), \quad \mathcal{P} = \mathcal{T}(P), \quad \mathcal{E} = \mathcal{T}(E),$$
$$\hat{\eta} = \mathcal{T}(\eta), \quad \hat{\nabla}_Y = \mathcal{T}(\nabla_Y), \quad \text{and} \quad \mathcal{Y} = \mathcal{Y}^* \text{ when not stated otherwise.}$$

Multi-Layer Perceptron

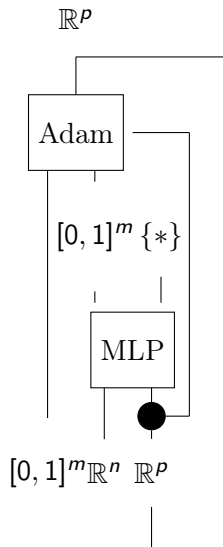
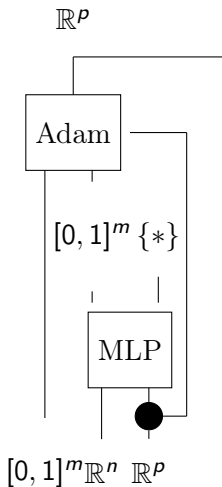
A classic multi-layer perceptron (MLP) classifier with an Adam optimizer is an instance of an abstract learning agent whose translator is defined as:

MLP

$$\begin{aligned}\mathcal{X} &= (\mathbb{R}^n)^{\mathbb{N}}, & \mathcal{Y} = \mathcal{Y}^* &= ([0, 1]^m)^{\mathbb{N}}, & \hat{\eta}_i &= \text{MLP}, \\ \mathcal{P} &= (\mathbb{R}^p)^{\mathbb{N}}, & \hat{\nabla}_{\mathcal{Y}_i} &= \text{Adam optimizer}, & \mathcal{E} &= \{*\}^{\mathbb{N}}.\end{aligned}$$



$\rightarrow \mathcal{T} \rightarrow$



Multi-Layer Perceptron

In this setting, the MLP is modeled by fixing each morphism component $\hat{\eta}_i = \text{MLP}$, independent of previous inputs. By removing this constraint, we can model broader classes of learning agents such as:

- Recurrent Neural Networks (RNNs),
- Hopfield Networks,
- Transformers.

Neural Architecture Search

A classical Neural Architecture Search (NAS) algorithm is an instance of an abstract learning agent whose translator functor is defined as:

NAS

$$\begin{aligned}\mathcal{X} &= (\mathbb{R}^n)^{\mathbb{N}}, & \mathcal{Y} &= \mathcal{Y}^* = ([0, 1]^m)^{\mathbb{N}}, & \hat{\eta}_i &= \text{MLP}_i, \\ \mathcal{P} &= (\mathbb{R}^p)^{\mathbb{N}}, & \hat{\nabla}_{\mathcal{Y}_i} &= \text{Adam optimizer}, & \mathcal{E} &= \{*\}^{\mathbb{N}}.\end{aligned}$$

Here, each MLP_i represents a different neural architecture at each step.

Explanation

Definition (Explanation)

Given an institution I , an object $\Sigma \in \text{Sign}_I$, and a concrete explainer $\hat{\eta} = \mathcal{T}(\eta) : \mathcal{X} \times \mathcal{P} \rightarrow \mathcal{Y} \times \mathcal{E}$, an *explanation* $\mathcal{E} = \mathcal{T}(E)$ in a language Σ is:

- a set of Σ -sentences (*syntactic explanation*), or
- a model of a set of Σ -sentences (*semantic explanation*).

Propositional Logic Explainer

Let I_{PL} be the institution of Propositional Logic and Σ a signature of I_{PL} such that

$$\{x_{flies}, x_{animal}, x_{plane}, x_{dark_color}, \dots\} \subseteq \Sigma,$$

with the standard connectives of Boolean Logic ($\neg, \wedge, \vee, \rightarrow$).

For instance, $\hat{\eta}$ could be an explainer aiming to predict an output in $\mathcal{Y} = \{x_{plane}, x_{bird}, \dots\}$ given an input in \mathcal{X} .

A **syntactic explanation** could be a Σ -sentence such as

$$\varepsilon = x_{flies} \wedge \neg x_{animal} \rightarrow x_{plane},$$

while a **semantic explanation** could be the truth-function of ε .

Isomorphism vs. Equivalence in XAI

Isomorphism (Local Equivalence)

$$\text{Mod}_{\Sigma_2}(\sigma(\varepsilon_1)) \cong \text{Mod}_{\Sigma_2}(\varepsilon_2)$$

Two explanations have identical meaning and structure.

Equivalence of Categories (Global Structural Equivalence)

$$\text{Mod}(\Sigma_1) \simeq \text{Mod}(\Sigma_2)$$

Whole explanation frameworks are structurally similar, translatable via functors preserving truth.

Translators for LA and XLA

We distinguish between two translators:

- \mathcal{T} with $\mathcal{T}(\mu) = \hat{\mu}$ for a LA
- \mathcal{T}' for a XLA

Objects of the latter are denoted with a prime:

$$\mathcal{T}'(Y) = \mathcal{Y}', \quad \mathcal{T}'(E) = \mathcal{E}'.$$

Post-hoc Explainer

Definition (Post-hoc Explainer)

Given a trained Learning Agent (LA) model

$$\hat{\mu} : \mathcal{X} \times \mathcal{P} \rightarrow \mathcal{Y},$$

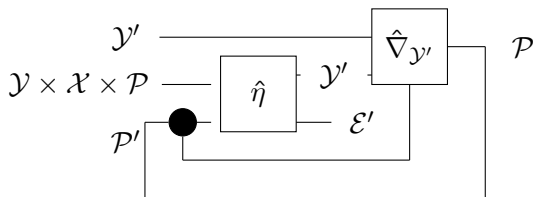
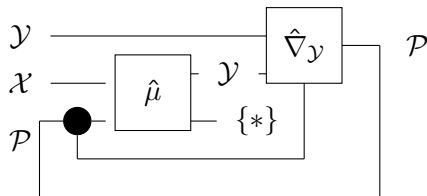
a *post-hoc explainer* is an Explaining Learning Agent (XLA) instantiated by a translator \mathcal{T}' such that

$$\hat{\eta} : \mathcal{X}' \times \mathcal{P}' \rightarrow \mathcal{Y}' \times \mathcal{E}', \quad \text{with } \mathcal{X}' = \mathcal{Y} \times \mathcal{X} \times \mathcal{P}.$$

Intuition:

- The explainer $\hat{\eta}$ operates on both inputs and outputs of $\hat{\mu}$.
- It produces an explanation \mathcal{E}' describing the behavior of $\hat{\mu}$.
- Feedback is mediated by $\nabla'_{\mathcal{Y}}$ in the XLA framework.

Post-hoc Explainer



Intrinsic Explainer

Definition (Intrinsic Explainer)

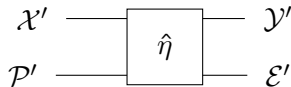
An *intrinsic explainer* is an Explaining Learning Agent (XLA) $\hat{\eta}$ whose input objects are parameters \mathcal{P}' and a set of entries of a database \mathcal{X}' :

$$\hat{\eta} : \mathcal{X}' \times \mathcal{P}' \rightarrow \mathcal{Y}' \times \mathcal{E}'.$$

Intuition:

- The explanation process is *intrinsic* to the model.
- Explanations \mathcal{E}' are generated directly during training or inference.
- No separate post-hoc model is required (interpretability is built into $\hat{\eta}$).

Intrinsic Explainer



Model-agnostic Explainer

Definition (Model-agnostic Explainer)

Given a Learning Agent (LA) model

$$\hat{\mu} : \mathcal{X} \times \mathcal{P} \rightarrow \mathcal{Y},$$

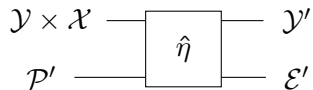
a *model-agnostic explainer* is an Explaining Learning Agent (XLA) $\hat{\eta}$ defined as

$$\hat{\eta} : \mathcal{X}' \times \mathcal{P}' \rightarrow \mathcal{Y}' \times \mathcal{E}', \quad \text{with } \mathcal{X}' = \mathcal{Y} \times \mathcal{X}.$$

Intuition:

- The explainer $\hat{\eta}$ operates independently of the internal structure of $\hat{\mu}$.
- It only requires access to inputs \mathcal{X} and outputs \mathcal{Y} .
- Examples include LIME and SHAP.

Model-agnostic Explainer



Model-specific Explainer

Definition (Model-specific Explainer)

Given a Learning Agent (LA) model

$$\hat{\mu} : \mathcal{X} \times \mathcal{P} \rightarrow \mathcal{Y},$$

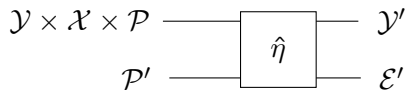
a *model-specific explainer* is an Explaining Learning Agent (XLA) $\hat{\eta}$ that depends on the internal structure or parameters of the model being explained:

$$\hat{\eta} : \mathcal{X}' \times \mathcal{P}' \rightarrow \mathcal{Y}' \times \mathcal{E}', \quad \text{with } \mathcal{X}' = \mathcal{Y} \times \mathcal{X} \times \mathcal{P}.$$

Intuition:

- The explainer $\hat{\eta}$ uses the model's parameters or gradients.
- It is *specific* to the structure of $\hat{\mu}$.
- Examples include Grad-CAM or Integrated Gradients.

Model-specific Explainer



Forward-based Explainer

Definition (Forward-based Explainer)

Given a gradient-based Learning Agent (LA) model

$$\hat{\mu} : \mathcal{X} \times \mathcal{P} \rightarrow \mathcal{Y},$$

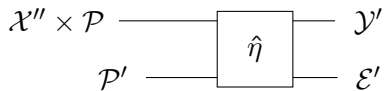
a *forward-based explainer* is an Explaining Learning Agent (XLA) $\hat{\eta}$ defined as

$$\hat{\eta} : \mathcal{X}' \times \mathcal{P}' \rightarrow \mathcal{Y}' \times \mathcal{E}', \quad \text{with } \mathcal{X}' = \mathcal{X}'' \times \mathcal{P}.$$

Intuition:

- Operates on the **forward pass** of the LA.
- Uses model activations and parameters to produce explanations.
- Examples: *Activation-based relevance methods*, feature attributions.

Forward-based Explainer



Backward-based Explainer

Definition (Backward-based Explainer)

Given a gradient-based Learning Agent (LA) model

$$\hat{\mu} : \mathcal{X} \times \mathcal{P} \rightarrow \mathcal{Y},$$

and an optimizer

$$\nabla_Y : \mathcal{Y} \times \mathcal{Y} \times \mathcal{P} \rightarrow \mathcal{P},$$

a *backward-based explainer* is an Explaining Learning Agent (XLA) $\hat{\eta}$ such that

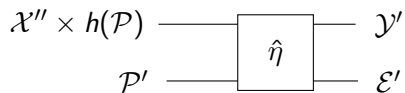
$$\hat{\eta} : \mathcal{X}' \times \mathcal{P}' \rightarrow \mathcal{Y}' \times \mathcal{E}', \quad \text{with } \mathcal{X}' = \mathcal{X}'' \times h(\mathcal{P}),$$

where $h(\mathcal{P}) = \frac{\partial \mathcal{L}(\mathcal{Y}, \mathcal{Y})}{\partial \mathcal{P}}$ is the gradient of the loss function.

Intuition:

- Operates on the **backward pass** (gradients) of the model.
- Uses gradient information to derive explanations.
- Examples: *Grad-CAM*, *Integrated Gradients*, *DeepLIFT*.

Backward-based Explainer



Generalized Learning Schemes

Learning schemes	$\mathcal{T}(X)$	$\mathcal{T}(Y)$	$\mathcal{T}(Y^*)$	$\mathcal{T}(E)$	$\mathcal{T}(\eta)$
Gen. unsup. model	$\diamond^{\mathbb{N}}$	$\diamond^{\mathbb{N}}$	$\{*\}^{\mathbb{N}}$	$\{*\}^{\mathbb{N}}$	\star
Gen. sup. model	$\diamond^{\mathbb{N}}$	$\diamond^{\mathbb{N}}$	$\diamond^{\mathbb{N}}$	$\{*\}^{\mathbb{N}}$	\star
Gen. continual lear. model	$\diamond^{\mathbb{N}}$	$\diamond^{\mathbb{N}}$	$\diamond^{\mathbb{N}}$	$\{*\}^{\mathbb{N}}$	\star
Gen. explaining model	$\diamond^{\mathbb{N}}$	$\diamond^{\mathbb{N}}$	$\diamond^{\mathbb{N}}$	$\diamond^{\mathbb{N}}$	\star

Resources

- Barbiero, P., Fioravanti, S., Giannini, F., Tonda, A., Lio, P., Di Lavore, E. (2023). Categorical Foundations of Explainable AI: A Unifying Theory. arXiv. [Link](#)
- Goguen, J. A., Burstall, R. M. (1992). Institutions: Abstract Model Theory for Specification and Programming. [Link](#)
- Hodges, W. (1997). A Shorter Model Theory. Cambridge University Press. [Link](#)
- Nakahira, K. (2023). Diagrammatic category theory. arXiv. [Link](#)
- Riehl, E. (2014). Category Theory in Context. (online lecture notes) [Link](#)