# Introduction to *In-Silico* learning
# Modelling Learning

**Abstract**

This tutorial provides an introduction to learning from a theoretical perspective: this preliminary introduction is therefore required to understand the incoming algorithms and models. The next tutorials will provide more concrete examples providing specific instances of these problems.
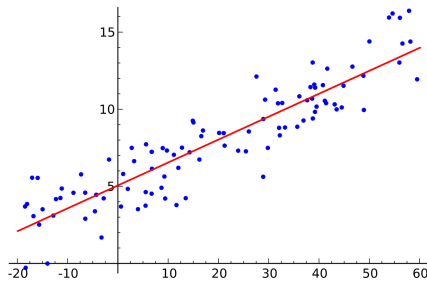
## 1  Intuition

In this tutorial we prefer using the term *In-Silico* learning to machine learning, because such term does not necessarily include other learning-based algorithms, such as data mining ones. In particular, both activities can be assimilated to a reasoning process, after which the machine can solve specific problems that would require "intelligence" for a human to solve. Therefore, these set of tutorials will use a different umbrella term for englobing two distinct approaches to automate learning [2].

Generally speaking, **machine learning** is the study of systems that improve their behaviour over time with experience: such experience consists in a set of examples $E$ in which we look for regularities or patterns fitting a specific model. In particular, these models focus on learning one specific function approximating the target function that the data implicitly exhibits. Figure 1a sketches a naïve learning approach: using a *linear least-squares (error) estimation*, we can infer a single hypothesis $h$ which, in this case, is a line minimizing the distance from the prediction (line's $y$ value) to the real data value (*loss*). Learning approaches then differ from the different model or function type that is going to be returned as an hypothesis.

Due to the limitations of the first attempts, more advanced symbolic reasoning approaches involving the use of both logic and probabilistic models emerged (*relational learning*[1]): in these approaches, as well as in **data mining** algorithms, we want to provide several different plausible hypotheses $h$ that meet some quality criterion $\mathcal{Q}$ over a sample $D \subseteq E$ of the original dataset. Similarly to the learning algorithms, $\mathcal{Q}$ and the choice of how to sample $D$ from $E$ might differ from algorithm to algorithm but, as we will see in the next section, they can be formalized with one single "mathematical"

---

[1]We're not going to cover these kind of algorithms now, as they require advance knowledge of probability theory and logic.



(a) Trying to fit the data (blue dots) to a linear model. The red line represents the outcome of the learning process.

| ID| | Items |
|---|---|
| 1 | {Bread, Milk} |
| 2 | {Bread, Diapers, Beer, Eggs} |
| 3 | {Milk, Diapers, Beer, Cola} |
| 4 | {Bread, Milk, Diapers, Beer} |
| 5 | {Bread, Milk, Diapers, Cola} |

| Frequent Item-sets | Association Rules |
|---|---|
| {Diapers, Beer} | {Diapers} → {Beer} |
| {Bread, Milk} | {Bread } → {Milk} |

(b) Using a frequency predicate $\mathcal{Q}$, we can determine the most frequent items over all the possible subset of the data items, thus providing multiple possible hypotheses.

Figure 1: Two examples of different approaches to *In-Silico* learning.

(a) Hypotheses viewed as functions.  (b) The *covers* relation for $\mathcal{Y} = 2$.
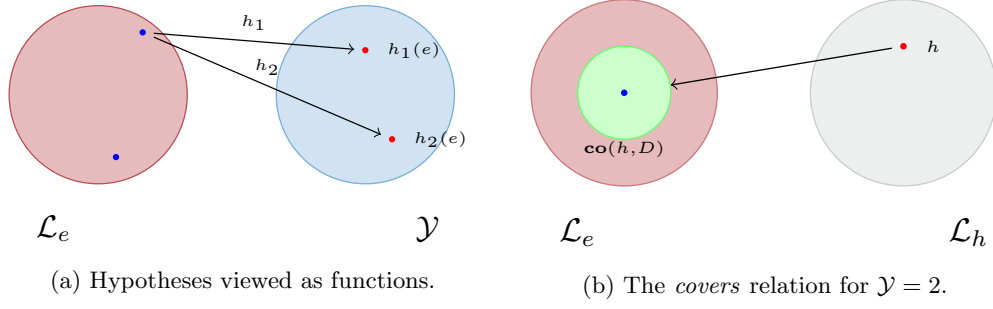
Figure 2: Examples of mathematical models that are used as intermediate steps for achieving computational representation.

formulation. An example of two data mining algorithms' outcome is provided in Figure 1b: using two distinct quality criterion over two distinct representations of the items, we can first detect the most frequent subsets, and then extract which are the most likely rules.

## 1.1 Formalization (*)

*We are now going to formalize such distinct approaches in order to outline differences and similarities of the two distinct approaches.*

All the records in our dataset $E$ come with a data **schema** $R(X_1, \ldots, X_n, Y)$. A schema is a table named $R$ with a set of *attributes* $X_1, \ldots, X_n, Y$, where the $X_i$ variables name the different "fields" of the training data, while $Y$ names the target classification name. Each **record** (or tuple, or data point) in $E$ is defined as a mapping (or function) $e$ of each attribute $Z \in \{X_1, \ldots, X_n, Y\}$ into a value in $dom(Z)$ [1]: the **training data** is represented by the set $\mathcal{L}_e \supseteq dom(X_1) \times \cdots \times dom(X_n)$, while the **target classification classes** are defined in $\mathcal{Y} \supseteq dom(Y)$.

Usually, each data point in $E$ is represented as a pair. $(\mathbf{x}, y) \in \mathcal{L}_e \times \mathcal{Y}$, so to clearly separate the training data from the classes value: from this definition, it follows that each pair-represented data point represents the map $(\mathbf{x}, f(\mathbf{x}))$ of a target function $f \colon \mathcal{L}_e \to \mathcal{Y}$ that we want to approximate with an hypothesis $h$ inferred from the data [2].

Please note that we are only interested to solve problems when $|\mathcal{Y}| \geq 2$, because for $|\mathcal{Y}| < 2$ the solution to the problem is always trivial. When $|\mathcal{Y}| = 2$, then we're going to solve a **binary classification problem**, while when $|\mathcal{Y}| > 2$ we must deal with a **multiclass classification problem**. For binary classification problems, we are going to use $\mathcal{Y} = \{0, 1\}$: all the examples such that $f(e) = 1$ are *positive*, while the others are called *negative*. This terminology comes from relational learning, where the former examples are the one that are logically true, while the latter are the ones to be considered as false.

**Example 1.** *Table 1 provides a dataset with the schema:*

$$\textsc{StarCraftReplay}(\textit{Age}, \textit{HoursPerWeek}, \ldots, \textit{MaxTimeStamp}, \textit{LeagueIndex})$$

*where we want to predict the XP level of the player (`LeagueIndex`) from all the remaining variables (`Age`,`HoursPerWeek`,…,`MaxTimeStamp`). We want to predict eight classes, and in fact $dom(\textit{LeagueIndex}) = \mathcal{Y} = \{1, 2, 3, 4, 5, 6, 7, 8\}$. Last, the first row $e$ can be represented either as a record with the fields mapping $e.\textit{Age} = 27 \ldots e.\textit{MaxTimeStamp} = 127448$, $e.\textit{LeagueIndex} = 5$ or as a pair $((27, 10, 3000, \ldots, 127448), 5)$. This means that the data induces a function $f(27, 10, 3000, \ldots, 127448) = 8$ that we need to learn using any approach of choice. We're going to use the former notation when accessing to a specific field is relevant (e.g., decision trees) and to the latter case in all the remaining situations.*

**Machine Learning** As per previous discussion, the aim of machine learning (or data mining) algorithm is to learn one (or multiple possible) **hypotheses**[2] $h \colon \mathcal{L}_e \to \mathcal{Y}$ assigning the training data to one single target classification class (Figure 2a). We can define the set of all the possible hypotheses

---

[2]Also called "decision function".

Table 1: A subset of the StarCraft II dataset available online at `http://summit.sfu.ca/item/13328`. This table subset provides the following fields: *(i)* the age of each player $(X_1, \ldots)$, *(ii)* the reported hours spent playing per week, *(iii)* the reported total hours spent playing, *(iv)* the number of continuous actions per minute, *(v)* the number of item selections via hotkeys per time frame, *(vi)* the number of assigned items via hotkeys per time frame, *(vii)* the number of unique and continuous hotkeys used per time frame, *(viii)* the number of continuous attack on a minimap per time frame, *(ix)* the number of continuous right clicks on a minimap per time frame, *(x)* the mean duration between PACs per time frame, *(xi)* the mean latency from the onset of PACs to their first action, *(xiii)* the mean number of actions within each PAC, *(xiii)* the total number of $24 \times 24$ game coordinate grids viewed by player per time frame, *(xv)* the number of battle items trained per time frame, *(xvi)* the number of unique units made per time frame, *(xvii)* the number of complex battle items trained per time frame, *(xviii)* the number of abilities requiring specific targeting instructions used per time frame, *(xix)* Time stamp of game's last recorded event $(\ldots, X_{19})$, and *(xx)* the target XP level class $(Y)$.

as[3] $\mathcal{Y}^{\mathcal{L}_e}$. As per previous discussions, the aim of machine learning algorithms is to find an hypothesis $h$ maximizing such loss function, which might be described as the discrepancy between the prediction by $h$ and the expected result. Given a distance[4] measure $\delta_{\mathcal{Y}}$ among the classes of interest, we can define the *least mean squares* loss function as follows:

$$loss_{\text{lms}}^{\delta_{\mathcal{Y}}}(h, E) := \sum_{(\mathbf{x}, y) \in E} \delta_{\mathcal{Y}}^2(y, h(\mathbf{x}))$$

For example, if $\mathcal{Y} \subseteq \mathbb{R}^n$ for any $n \in \mathbb{N}$, then we can use $\delta_{\mathbb{R}^n}^2$ as the euclidean distance $\delta_{\mathbb{R}^n}^2(\mathbf{x}, \mathbf{x}') := \|\mathbf{x} - \mathbf{x}'\|$. The training problem can be summarized as follows:

$$ML(loss, D, \mathcal{L}_h) := \arg\min_{h' \in \mathcal{L}_h} loss(h', E)$$

With respect to binary classification problems, in most real world algorithms, we're not necessarily going to learn a discrete $h\colon \mathcal{L}_e \to \{0, 1\}$, but we're going to learn a function approximating such values in $\mathbb{R}$ (e.g., a regression function), so $\tilde{h}\colon \mathcal{L}_e \to [0, 1]$ is learned instead. In these situations, we might define the hypothesis from the regression function as follows:

$$h(\mathbf{x}) := \arg\min_{y \in \mathcal{Y}} \delta_{\mathcal{Y}}(y, \tilde{h}(\mathbf{x}))$$

**Data Mining**  Let us now suppose that we deal with data mining problems, where a possible outcome of the algorithm might be a set of elements $B \subseteq \mathcal{L}_e$. We know from set theory that we can always express such set as an *indicator function* $\mathbf{1}_B\colon \mathcal{L}_e \to \{0, 1\}$ determining whether any element of $\mathcal{L}_e$ is also in $B$ or not[5]. In this scenario we can still model our hypotheses as per previous discussions: for each possible returned set $B$, we can return an hypothesis $h_B$ which is actually the indicator function $\mathbf{1}_B$. This also implies that each data mining algorithm will retrieve only the *predicted positive* examples.

Contrarily to machine learning algorithms, data mining algorithms are not going to return a single hypothesis, but all such hypotheses that satisfy a quality criterion $\mathcal{Q}$, that often requires $B$: when such set is not explicitly given but $h$ is returned instead, we can define a cover function $\mathbf{co}(h, D) = \{x \in D \mid h(x) = 1\}$ identifying the set of all the items that are covered by the hypothesis $h$. An intuitive graphical representation of such function is given at Figure 2b. Data mining algorithms are often interested in returning only the hypothesis exhibiting frequent patterns from the data: in order to do so, we can define the *min-frequency* predicate that is going to return only the hypothesis above a given threshold $\theta$:

$$\mathcal{Q}_\theta(h, D) := |\mathbf{co}(h, D)| \geq \theta$$

This means that we can generalize all the possible data mining problems as follows:

$$DM(\mathcal{Q}, D, \mathcal{L}_h) := \{h \in \mathcal{L}_h \mid \mathcal{Q}(h, D) \text{ holds}\}$$

Albeit it might seem quite inefficient to visit all the possible hypotheses in $\mathcal{L}_h$, we might provide an intuition[6] of how we can efficiently do so by pruning most of the non-valid hypotheses: for binary classification problems, the set of all the possible hypotheses is $2^{\mathcal{L}_e}$, so each possible predicted positive example in $\wp(\mathcal{L}_e)$ represents the outcome of one single hypothesis. Given that there is an isomorphism between $\mathcal{L}_h$ and $\wp(\mathcal{L}_e)$ (Lemma A), we can search directly the hypotheses from the latter space. Let us now also suppose that we want to use the min-frequency predicate as a possible quality criterion: given that all the generalizations of an hypothesis will satisfy such criterion (Lemma 2), then as a dual condition none of the specializations of the hypotheses not-satisfying such criterion will satisfy it [2]. As a result, we can define top-down search algorithms starting from the most general hypothesis and specializing it: we can stop expanding an hypothesis when such hypothesis doesn't meet the quality criterion.

---

[3]See `https://en.wikipedia.org/wiki/Power_set#Representing_subsets_as_functions`.

[4]See the first tutorial on algorithms and data structure for a formal definition of a distance metric.

[5]$\mathbf{1}_B(\mathbf{x}) := \begin{cases} 1 & \mathbf{x} \in B \\ 0 & \mathbf{x} \notin B \end{cases}$

[6]For the formal proofs see Appendix A.

**Prediction outcome**

|  |  | p | n | total |
|---|---|---|---|---|
| **actual value** | **p′** | True Positive | False Negative | P′ |
|  | **n′** | False Positive | True Negative | N′ |
| **total** |  | P | N |  |

Figure 3: Displaying the performance of a learning algorithm.

# 2 Multiclass Classification from Binary Classification

At this stage, we mainly discussed binary classification problems. This decision is based on the fact that both (1) these classifiers are simple enough to be understood in a few days time, and (2) it is always possible to reduce an arbitrary classification problem to a binary classification problem.

- **one-versus-all**: for each class $y \in \mathcal{Y}$, train a binary (e.g. learning) algorithm so that it will learn whether any $x \in \mathcal{L}_h$ belongs to the class $y$ or not. For each of these algorithms, obtain the approximated hypothesis function $\tilde{h}_y$. Now, we can define the ensemble hypothesis as follows:

$$h_{\mathrm{OvA}}(\mathbf{x}) = \arg\max_{y \in \mathcal{Y}} \tilde{h}_y(\mathbf{x})$$

- **one-versus-one**: for each class $y \in \mathcal{Y}$, initialize a binary (e.g., learning) algorithm $\mathcal{A}_{(y_1, y_2)}$ for each distinct pair $(y_1, y_2) \in \mathcal{Y}^2$, infer the binary hypothesis $\mathrm{Im}(h_{(y_1, y_2)}) = \{ y_1, y_2 \}$ and store it in a pool ($|C| = \frac{|\mathcal{Y}|(|\mathcal{Y}|-1)}{2}$ ). Then, associate $\mathbf{x}$ to the class winning the following maximum vote strategy.

$$h_{\mathrm{OvO}}(\mathbf{x}) = \arg\max_{y \in \mathcal{Y}} | \{ h_{(y_1, y_2)} \in C \mid h_{(y_1, y_2)}(\mathbf{x}) = y \} |$$

Furthermore, this technique is independent from the specific in-silico model of choice, given that it only deals with the returned hypotheses.

# 3 Confusion Matrix for Binary Classification

Last, we can use confusion matrices to display the performance of a classification algorithm. In fact, these matrices can be used for both determine how many predicted positive elements are actually true (**precision**), and how many positive elements were actually retrieved even though wrongly classified (**recall**). Based on the matrix in Figure 3, we can define precision and recall as the following metrics:

$$\mathbf{Precision} = \frac{|p \cap p'|}{|P|} \qquad \mathbf{Recall} = \frac{|p \cap p'|}{|P'|}$$

This confusion matrix can be easily generalized for the multi-class classification problem.

# References

[1] Christopher M. Bishop. *Fundamentals of Database Systems.* Pearson, 7th edition, 2015.

[2] Luc De Raedt. *Logical and Relational Learning.* Springer Verlag, Germany, 2008.

# A   Formal Proofs

Let us consider binary classification problems. One natural way to structure the search space is to employ the generality relation $\preceq$: we can state that $h_2$ is more general that $h_1$ if all the examples covered by $h_1$ are also covered by $h_2$, that is:

$$h_2 \preceq h_1 \Leftrightarrow \mathbf{co}(h_1, D) \subseteq \mathbf{co}(h_2, D)$$

Now, given that the hypotheses are just indicator functions for sets $B_i \in \wp(\mathcal{L}_e)$, we can express such relations in terms of the predicted positive sets

**Lemma 1.** *Under the assumption that $h_1 = \mathbf{1}_{B_1}$ and $h_2 = \mathbf{1}_{B_2}$, then for binary decision problems we have that $h_2 \preceq h_1 \Leftrightarrow B_1 \subseteq B_2$.*

*Proof.* If we assume to know that, for any subset $S \subseteq T$ of $T$, the intersection $S \cap T$ is equal to $S$[7] (**T**), then we can provide the following proof:

$$
\begin{aligned}
h_2 \preceq h_1 &\Leftrightarrow \mathbf{co}(h_1, D) \subseteq \mathbf{co}(h_2, D) \\
&\Leftrightarrow (B_1 \cap D) \subseteq (B_2 \cap D) \\
&\overset{\mathbf{T}}{\Leftrightarrow} (B_1 \cap D) \cap (B_2 \cap D) = (B_1 \cap D) \\
&\Leftrightarrow (B_1 \cap B_2) \cap D = (B_1) \cap D \\
&\Leftrightarrow B_1 \cap B_2 = B_1 \\
&\Leftrightarrow B_2 \subseteq B_2
\end{aligned}
$$

$\square$

We want now to show that min-support can be used as a pruning quality criterion while visiting $\wp(\mathcal{L}_e)$. First, we need to prove that such criterion is anti-monotonic. The definition of anti-monotonicity for $h_1, h_2$ is given as follows:

$$\forall D \subseteq \mathcal{L}_e. \ (h_1 \preceq h_2) \wedge \mathcal{Q}(h_2, D) \Rightarrow \mathcal{Q}(h_1, D)$$

**Lemma 2.** *The min-support quality criterion is anti-monotonic with respect to a given set $D$.*

*Proof.* By Lemma 1, we can rewrite $h_1 \preceq h_2$ as $B_2 \subseteq B_1$ (**H1**). Also, the min-support quality criterion for $h_2$ becomes $|B_2 \cap D| \geq \theta$, thus implying that we know a $x \in \mathbb{R}^+$ for which $|B_2 \cap D| = \theta + x$ (**H2**). If we now expand the min-frequency definition for $h_1$, then we have to prove that $|B_1 \cap D| \geq \theta$. By expanding **H1**, we have that we know a set $S$ for which $|(B_2 \cup S) \cap D| = |(B_2 \cap D) \cup (S \cap D)| \geq \theta$. We can now prove:

$$
\begin{aligned}
|(B_2 \cup S) \cap D| = |(B_2 \cap D) \cup (S \cap D)| \geq \theta &\Leftrightarrow \exists y \in \mathbb{R}^+. |B_2 \cap D| + y = \theta \\
&\overset{\mathbf{H2}}{\Leftrightarrow} \exists y \in \mathbb{R}^+. \theta + x + y = \theta \\
&\Leftrightarrow \theta \geq \theta
\end{aligned}
$$

$\square$

---

[7]See `https://proofwiki.org/wiki/Intersection_with_Subset_is_Subset`