

(c) Filip Železný, all rights reserved. Any use of this material only with prior approval by the author.

# Symbolic Machine Learning

Filip Železný

## Contents

# 1 A General Framework

## 1.1 Percepts and Actions

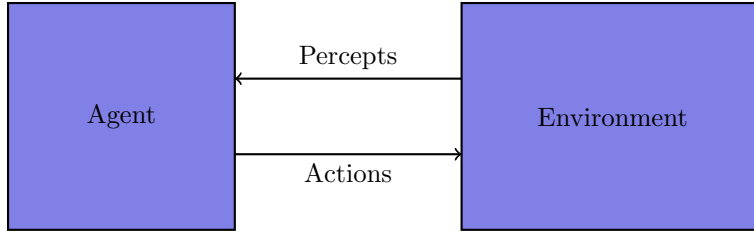


Figure 1: The basic situation under study.

- Discrete *time*

$$k = 1, 2, \dots$$

- *Percepts*

$$\forall k : x_k \in X$$

- *Actions*

$$\forall k : y_k \in Y$$

$X$  and  $Y$  are finite.

A *history* is a sequence of alternating percepts and actions, i.e.,

$$x_1, y_1, x_2, y_2, \dots, x_k, y_k$$

and is denoted as  $xy_{\leq k}$ . Similarly,  $xy_{< k} = x_1, y_1, x_2, y_2, \dots, x_{k-1}, y_{k-1}$ . There is a probability distribution  $\mu$  on histories

$$\mu(xy_{\leq k}) = \mu(x_1)\mu(y_1|x_1)\mu(x_2|x_1, y_1) \dots \mu(x_k|xy_{< k})\mu(y_k|x_k, xy_{< k}) \quad (1)$$

After the initial ‘kick-off’  $x_1$  from the environment distributed according to  $\mu(x_1)$ , any percept  $x_k$  generated by the environment at time  $k$  depends on the entire preceding history  $xy_{< k}$  according to

$$\mu(x_k|xy_{< k}) \quad (2)$$

Actions  $y_k$  are determined by agent’s decision *policy* which also depends on the history as well as the current percept and are distributed according to

$\mu(y_k|x_k, xy_{<k})$ . We will assume that the policy is *deterministic*. Thus we identify the policy with function  $\pi : (X \times Y)^* \times X \rightarrow Y$ , so

$$y_k = \pi(xy_{<k}, x_k) \tag{3}$$

This means that  $\mu(y_k|xy_{<k}, x_k) = 1$  for  $y_k = \pi(xy_{<k}, x_k)$  and 0 otherwise.

The following diagram illustrates the influences between the introduced variables.

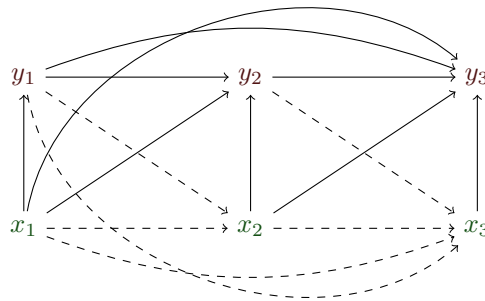


Figure 2: Influence diagram for actions  $y_k$  and percepts  $x_k$  for  $1 \leq k \leq 3$  with full lines indicating deterministic influences (via  $\pi$ ) and dashed lines showing probabilistic influences (via  $\mu$ ).

While we have yet to define what goals the agent should achieve through interaction with the environment, obviously some histories will be “better” than others in terms of the goal achievement. To maximize the probability (1) of good histories, the agent cannot influence the conditional probability (2), which is inherent to the environment, but it can follow a good policy (3). However, the effect of actions proposed by the policy depends on (2) which is generally not known to the agent. So the agent needs to recognize the environment by experimenting with it. This is formally reflected by (3) where action  $y_k$  depends not only on the current percept  $x_k$  but also on the history  $xy_{<k}$ . So the agent will generally make different decisions  $y_k \neq y_{k'}$  for  $k > k'$  even if  $x_k = x_{k'}$  because the experience  $xy_{<k}$  at time  $k$  is larger than experience  $xy_{<k'}$  at time  $k'$ . This is our first reflection of *learning*.

How does the agent know how well it is doing? This information comes from the environment through a specially distinguished part of the percepts, called *rewards*. The remaining part of each percept contains *observations*. Formally,  $X = O \times R$ ,  $o_k \in O$ ,  $r_k \in R \subset \mathfrak{R}$ , so

$$x_k = (o_k, r_k) \tag{4}$$

Since  $X$  is assumed finite, it follows that rewards have their finite minimum and maximum.

The probability of  $x_k$  in (2) can be written in terms of the marginals  $\mu_O$  and  $\mu_R$

$$\begin{aligned}\mu(x_k|xy_{<k}) &= \mu(o_k, r_k|xy_{<k}) = \\ \mu_O(o_k|r_k, xy_{<k})\mu_R(r_k|xy_{<k}) &= \mu_R(r_k|o_k, xy_{<k})\mu_O(o_k|xy_{<k})\end{aligned}$$

which also makes it clear that  $o_k$  and  $r_k$  are in general not assumed mutually independent, even if conditioned on  $xy_{<k}$ .

## 1.2 Nonsequential Cases

Scenarios where current percepts depend on the history of previous percepts and actions are called *sequential*. The framework described so far is maximally general in that dependence is assumed on the entire history from  $k = 1$  on. On the other extreme are *nonsequential* scenarios. Here, observations are independent of the history as well as the current reward, i.e.

$$\mu_O(o_k|r_k, xy_{<k}) = \mu_O(o_k) \tag{5}$$

and thus  $o_1, o_2, \dots$  are mutually independent random variables sampled from the same distribution  $\mu_O$  (they are “i.i.d.”).

Rewards in the nonsequential case are assumed to depend only the immediately preceding observation and the action taken on it, i.e.

$$\mu_R(r_k|o_k, xy_{<k}) = \mu_R(r_k|o_{k-1}, y_{k-1}) \tag{6}$$

however, since  $y_{k-1}$  is functionally determined by the history  $xy_{<k-1}$  and percept  $x_{k-1} = (o_{k-1}, r_{k-1})$  through (3), we may rewrite (6) as

$$\mu_R(r_k|o_{k-1}, r_{k-1}, xy_{<k-1}) \tag{7}$$

which makes it clear that reward  $r_k$  depends on previous rewards, and thus rewards  $r_1, r_2, \dots$  are not i.i.d.. This is natural since if they were, it would mean the agent never improves its performance.

## 1.3 Batch Learning

We will also consider a specific yet important nonsequential case called *batch learning* consisting of two phases switching right after time  $K$

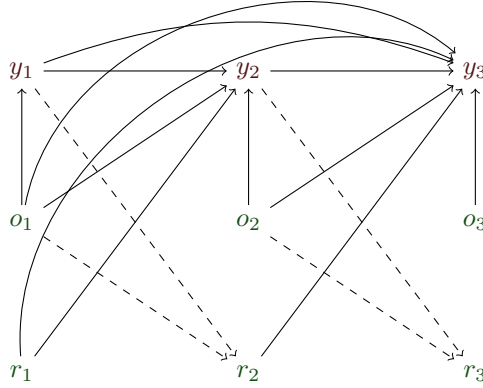


Figure 3: Influence diagram for actions  $y_k$ , observations  $o_k$ , and rewards  $r_k$  for  $1 \leq k \leq 3$  with full lines indicating deterministic influences (via  $\pi$ ) and dashed lines showing probabilistic influences (via  $\mu$ ) in the nonsequential case.

- the *learning (training, exploration) phase* at  $k = 1, 2, \dots, K$
- the *action (testing, exploitation) phase* taking place in  $k = K+1, K+2, \dots$

In the action phase, the agent no longer changes its decision making policy, so

$$\text{if } k, k' > K \text{ and } x_k = x_{k'} \text{ then } y_k = y_{k'} \quad (8)$$

and ignores rewards. So the action proposed by the policy depends only on the current observation and the history only up to time  $K$ . So for  $k > K$ , (3) changes here into

$$y_k = \pi(xy_{\leq K}, o_k) \quad (9)$$

and (6, 7) change into

$$\mu_R(r_k | o_{k-1}, y_{k-1}) = \mu_R(r_k | o_{k-1}, xy_{\leq K}) \quad (10)$$

because due to (9),  $y_{k-1}$  is determined by  $o_{k-1}$  and  $xy_{\leq K}$ . The observation  $o_{k-1}$  does not depend on rewards due to (5). So reward  $r_k$  does not depend on previous rewards  $r_{k'}, k > k' > K$ . Another way to say this is that rewards in the action phase are conditionally independent of each other, given the learning phase history:

$$\mu_R(r_k, r_{k'} | xy_{<K}) = \mu_R(r_k | xy_{<K}) \mu_R(r_{k'} | xy_{<K}) \quad (11)$$

The following figure illustrates the batch-learning situation.

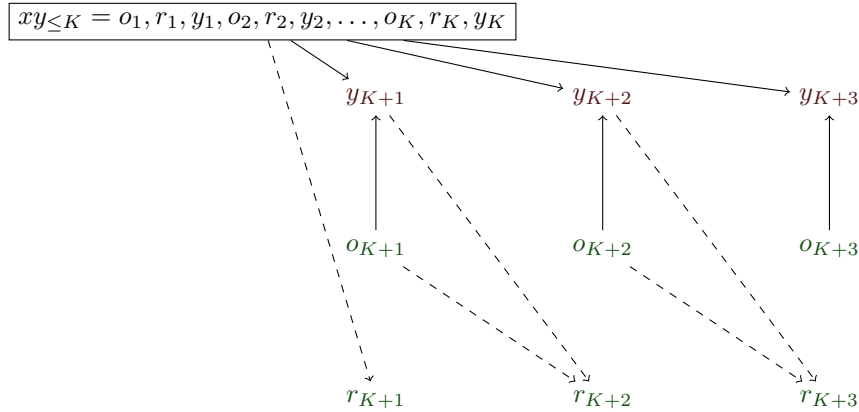


Figure 4: Influence diagram for actions  $y_k$ , observations  $o_k$ , and rewards  $r_k$  in the action phase ( $k > K$ ) of batch learning with full lines indicating deterministic influences (via  $\pi$ ) and dashed lines showing probabilistic influences (via  $\mu$ ). The top row indicates the influence of the learning phase on the agent’s decisions in the action phase.

We have observed that rewards  $r_k, r_{k'}$  (for any  $k, k' > K$ ) are mutually independent given the learning phase  $xy_{<K}$ . Note that they are also sampled from the same distribution. This may seem to contradict (10) which stipulates that  $r_k$  depends on the observation  $o_{k-1}$  whereas  $r_{k'}$  depends on  $o_{k'-1}$ . However, by (5),  $o_{k-1}$  and  $o_{k'-1}$  are sampled from the same distribution  $\mu_O$ . So we can express the distribution of  $r_k$  ( $\forall k > K$ ) without conditioning on the observations by marginalizing them away from the equation as follows

$$\mu_R(r_k | xy_{\leq K}) = \sum_{o_{k-1} \in \mathcal{O}} \mu_O(o_{k-1}) \mu_R(r_k | o_{k-1}, xy_{\leq K}) \quad (12)$$

So rewards in the action phase indeed are i.i.d. according to the above distribution conditioned only on the history of the learning phase.

## 1.4 Rewards and Goals

It has been obvious that the agent’s goal is to maximize rewards. Here we formalize this goal. Since rewards come at each point of the history, we want the agent to maximize their sum up to a finite time *horizon*  $m \in \mathbb{N}$

$$r_1 + r_2 + \dots + r_m$$

or, more generally, maximize the *discounted* sum

$$\sum_{k=1}^{\infty} r_k \gamma_k$$

where  $\forall k : \gamma_k \geq 0$  and  $\sum_{i=1}^{\infty} \gamma_i < \infty$ . The latter condition guarantees that the sum above converges, which is because the  $r_k$ 's are bounded by a constant (c.f. Section 1.1).

But since rewards are probabilistic, the agent should choose a sequence  $y_{\leq m}$  of actions leading to a high *expected* cumulative reward

$$\sum_{r_{\leq m}} \mu_R(r_{\leq m} | y_{\leq m}) (r_1 + r_2 + \dots + r_m)$$

or, in the discounted case

$$\lim_{m \rightarrow \infty} \sum_{r_{\leq m}} \mu_R(r_{\leq m} | y_{\leq m}) \sum_{k=1}^m r_k \gamma_k$$

where the first sum in both cases goes over all possible reward sequences  $r_{\leq m}$  (since  $R$  and  $m$  are finite, there is a finite number of them).

However, for the specific case of *batch learning*, we establish a more appropriate learning goal. First, we do not care about maximizing rewards in the *learning phase* as the purpose of this phase is to probe the environment even at the price of possibly poor rewards. Second, in the *action phase* after time  $K$ , the rewards  $r_k$ ,  $k > K$  are sampled independently from the same distribution (12) so we can simply maximize their expectation with respect to this distribution

$$\sum_{r_k \in R} \mu_R(r_k | xy_{\leq K}) r_k \tag{13}$$

It is again obvious from the formula that the expected reward only depends on the learning phase history  $xy_{\leq K}$ , after which the agent no longer changes its action policy. Note also that the batch learning scenario allowed us to define an objective (13) without the need to choose the parameters  $m$  or  $\gamma_k$  ( $k = 1, 2, \dots$ ) needed in the sequential scenario.

## 1.5 Environment States

With the exception of the non-sequential scenario, our framework has been very general in that percepts  $x_k$  generally depend on entire histories  $xy_{<k}$ . In the real world, many histories may be equivalent, i.e. leading to the same probabilities of  $x_k$  conditioned on action  $y_{k-1}$ . This can be formalized through the notion



of *environment state*. Intuitively, the state acts as the environment’s ‘memory’ carrying all the information from the history, which is important for generating percepts. To formalize this, we will assume there is a set  $E$  of all possible states, and instead of 2, we will prescribe that percepts at time  $k$  only depend on the state  $e_k \in E$  at that time, i.e. they are generated according to

$$\mu(x_k|e_k) \tag{14}$$

But how are the states determined? We will first explore a principle, which—in a sense—will turn out to be maximally general. In particular, assume that the initial state  $e_1$  is fixed to an ‘empty’ (or ‘dummy’) value  $e_1 = s^*$  such that

$$\mu(x_1|s^*) = \mu(x_1) \tag{15}$$

so the first percept is generated just as in (1). Afterwards, any state  $e_k$  ( $k > 1$ ) is established probabilistically by the preceding state, the last percept, and the last agent’s action through the following state *update* distribution

$$\mathcal{E}(e_k|e_{k-1}, x_{k-1}, y_{k-1}) \tag{16}$$

We said earlier that this principle was ‘maximally general.’ To say this more precisely, for any environment generating percepts by (2), we can set up  $E$  and  $\mathcal{E}$  such that

$$\mu(x_k|e_k) = \mu(x_k|xy_{<k}) \tag{17}$$

Indeed, if we allow  $E$  to be infinite, then there could simply exist a distinct state for each possible history (there is an infinite number of possible histories for unbounded  $k$ ). Then we can instantiate the distribution (16) to the functional dependence

$$e_k = e_{k-1} \parallel (x_{k-1}, y_{k-1}) \tag{18}$$

where  $\parallel$  denotes concatenation. As a result,  $e_k$  will simply collect the entire history and in (17),  $e_k$  would be just a different name for  $xy_{<k}$ .

However, we will make an important assumption, which will significantly simplify the framework, that the number of possible states is bounded by a finite constant  $E_{\max} \in N$  which does not depend on  $k$

$$|E| < E_{\max} \tag{19}$$

In practical tasks, there will be far fewer states than possible histories.

Moreover, we can afford an additional simplifying assumption which will further lessen the generality of the framework, while keeping it able to encompass the learning scenarios we are going to elaborate. In particular, we will assume quite naturally that the influence between environment states and the emitted percepts are single-directional in the sense that the percepts depend on states by (14) but not vice versa, so we remove  $x_{k-1}$  from (16)

$$\mathcal{E}(e_k|e_{k-1}, x_{k-1}, y_{k-1}) = \mathcal{E}(e_k|e_{k-1}, y_{k-1}) \tag{20}$$

As we have discussed already, the finiteness of  $E$  means that the environment has a ‘finite memory.’ Thus, from a current state  $e$ , one cannot in general identify the entire preceding history of states and actions. However, the framework as defined can still model environments that remember a *bounded number* of previous states and actions. We will demonstrate this through an example with a one-step memory.

Consider an environment with states  $E$  and update distribution  $\mathcal{E}$ , and assume that the current state  $e_k$  does not allow to infer the previous state  $e_{k-1}$  or action  $y_{k-1}$ . Then we can always define an extended (yet still finite) set of states as

$$E_{\text{ext}} = E \times E \times Y \quad (21)$$

(with the latter two factors acting as memory) and an extended update distribution

$$\mathcal{E}_{\text{ext}} \{(e_k, \text{olde}_k, \text{old}y_k) \mid (e_{k-1}, \text{olde}_{k-1}, \text{old}y_{k-1}), y_{k-1}\} \quad (22)$$

such that

- $e_k$  is distributed according to  $\mathcal{E}(e_k \mid e_{k-1}, y_{k-1})$
- $\text{olde}_k = e_{k-1}$  and  $\text{old}y_k = \text{old}y_{k-1}$ , both with probability 1.

So the three components of the extended state correspond, respectively, to the original state at current time  $k$ , the same at the previous time  $k - 1$  and the agent’s actions at time  $k - 1$ .

We will often model environments with a natural (interpretable) set of states  $E$ , producing percepts that depend not only on the current state but also on the state and agent’s action one-step back in history. We have just seen that such a situation can still be modeled with the simple assumption (14) by extending  $E$  towards a state set with a memory as in the above example. However, this leads to some cumbersome notation as in (22). We will avoid these complications by keeping the original state set  $E$  rather than extending it with memory for  $e_{k-1}$  and  $y_{k-1}$ . Instead, we will add the latter two explicitly to the conditional part of (14), thus obtaining

$$\mu(x_k \mid e_k, e_{k-1}, y_{k-1}) \quad (23)$$

Regarding the two components of the percepts  $x_k$ , the observations will depend only on the current state (and not on the previous one) and the last agent’s action

$$\mu(o_k \mid e_k, e_{k-1}, y_{k-1}) = \mu_o(o_k \mid e_k, y_{k-1}) \quad (24)$$

and the reward will depend on the previous state (and not on the current one) and the action taken immediately on it

$$\mu(r_k \mid e_k, e_{k-1}, y_{k-1}) = \mu_r(r_k \mid e_{k-1}, y_{k-1}) \quad (25)$$

The formulation (23)–(25) is convenient for interpretability and to avoid the complex notation such as in (22). One should however keep in mind that with a suitable definition of the state variable, it is possible to avoid the variables  $e_{k-1}, y_{k-1}$  in the conditional part and adhere to the simple prescription (14).

## 1.6 Agent States

A reasoning similar to the previous section applies to the agent, whose actions generally depend on the entire history as in (3). Again, many histories can lead to the same mapping from percepts to actions, for example because the agent has built the same hypothesis about the environment throughout the different histories. So analogically to the environmental states, we introduce the notion of agent’s *state*  $a \in A$  and postulate that

$$|A| < A_{\max} \tag{26}$$

for some constant  $A_{\max} \in \mathbb{N}$ .

We adopt a counterpart of (14), meaning that the action will be determined by the agent’s state, again through a functional prescription

$$y_k = \pi(a_k) \tag{27}$$

The current action thus does not depend explicitly on the current percept  $x_k$ . This is because the latter can be simply stored as a part of the state equipped with memory as can be shown through a reasoning very similar to that in the previous section. However, the dependence on the current percept, and specifically on its observation part, will be so typical that we will make it explicit. Again, this will save us notational complications entailed by the need to memorize percepts within states. We will thus use the formula

$$y_k = \pi(a_k, o_k) \tag{28}$$

Regarding the update rule, analogically to (20) we will assume that the state is updated (deterministically) given the previous one and the current percept.

$$a_k = \mathcal{A}(a_{k-1}, x_k) \tag{29}$$

This seems not fully analogical to (20) as the current percept, rather than the previous one is taken into account. This is due to our setting of the agent-environment communication, in which  $y_{k-1}$  is the last action taken before the environment updates its state, whereas  $x_k$  is the last percept received by the agent for its state update.

The formalization using environment states and agent hypotheses results in the agent and environment structures depicted in Fig. 5. The diagram of variable influences is shown in Fig. 6.

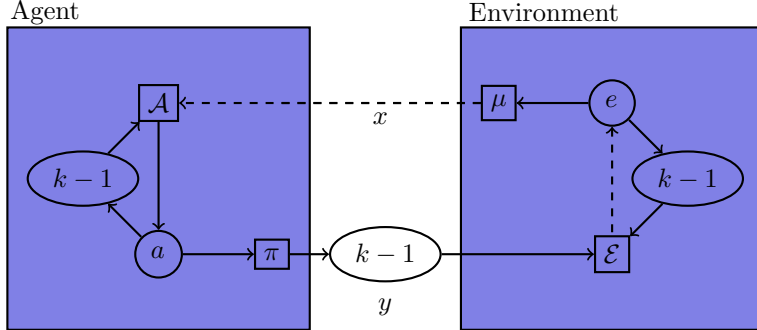


Figure 5: The state-based scheme of agent-environment interaction according to (14), (20) and (27), and (29). Full and dashed lines denote functional and probabilistic influences, respectively. The  $k - 1$  nodes denote a one-step time lag. Note that we have introduced further dependencies through (23) and (28), which are not shown in the picture. These dependencies can be avoided through a reformulation of the state variable and the state update function.

## 1.7 Nonsequential and Batch Cases with States and Hypotheses

Just like in the framework using entire histories, also with the formulation based on states and hypotheses the situation simplifies a lot in the *nonsequential case*. Here, the environment has no memory at all so the conditioning factors in (20) and states are updated by i.i.d. sampling from the marginal distribution

$$\mathcal{E}(e_k) \quad (30)$$

Furthermore, observations  $o_k$  no longer depend on agent's last action as in (24) so they are sampled from

$$\mu_o(o_k|e_k) \quad (31)$$

Since  $e_k$ 's are i.i.d., the  $o_k$ 's are also i.i.d.

Rewards, given by (25), are however still generally non-i.i.d. as they depend on the agent's actions, which in turn depend on the evolving agent's hypothesis. Fig. 7 shows the complete set of influences in the nonsequential case.

A further simplification comes in the special *batch-learning* scenario of the non-sequential case. While in the learning phase of the latter, the agent uses the update rule (29), in the action phase it no longer updates its state, so

$$a_k = a_K, \forall k \geq K \quad (32)$$

This is illustrated in Fig. 8. Special attention is needed regarding the variables at time  $K$ . Reward  $r_K$  (part of percept  $x_K$ ) is the last training reward, according

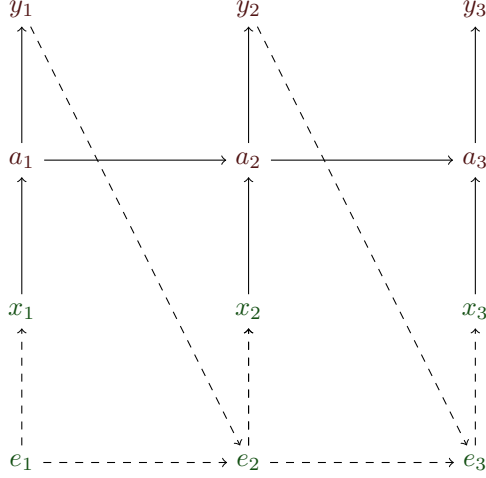


Figure 6: Influence diagram for states  $a$ , actions  $y_k$  and percepts  $x_k$  for  $1 \leq k \leq 3$  with full lines indicating deterministic influences (via  $\pi$  and  $\mathcal{A}$ ) and dashed lines showing probabilistic influences (via  $\mu$  and  $\mathcal{E}$ ). See caption to Fig. 5 for further relevant remarks.

to which the last update is conducted towards the final  $h_K$ . Observation  $o_K$  (another part of percept  $x_K$ ) is, however, the first *testing* observation.

Since  $y_{k-1}$  is determined by  $o_{k-1}$  and  $a_{k-1}$  through (28), we can rewrite (25) as

$$\mu_r(r_k | e_{k-1}, o_{k-1}, a_{k-1}) \quad (33)$$

Because environment states  $e_k$  and observations  $o_k$  in the non-sequential setting are i.i.d. (sampled from the respective distributions (30) and (31) independent of  $k$ ), we can marginalize the above as

$$\mu_r(r_k | a_{k-1}) = \sum_{e_{k-1} \in E} \sum_{o_{k-1} \in O} \mu_r(r_k | a_K, e_{k-1}, o_{k-1}) \mu_O(o_{k-1} | e_{k-1}) \mathcal{E}(e_{k-1}) \quad (34)$$

In the testing phase ( $k > K$ ) where  $a_k$  is fixed to  $a_k = a_K$ , rewards  $r_k$  are thus i.i.d. according to the distribution  $\mu_r(r_k | a_K)$  depending only on the final state  $a_K$  of training. This is analogical to the state-free formulation (12). Similarly to (13), an agent operating in the batch-learning scenario with states will be assessed by the expected reward in the testing phase

$$\sum_{r_k \in R} \mu_R(r_k | a_K) r_k \quad (35)$$

so in the training phase, it should reach a state  $a_K$  maximizing this quantity.

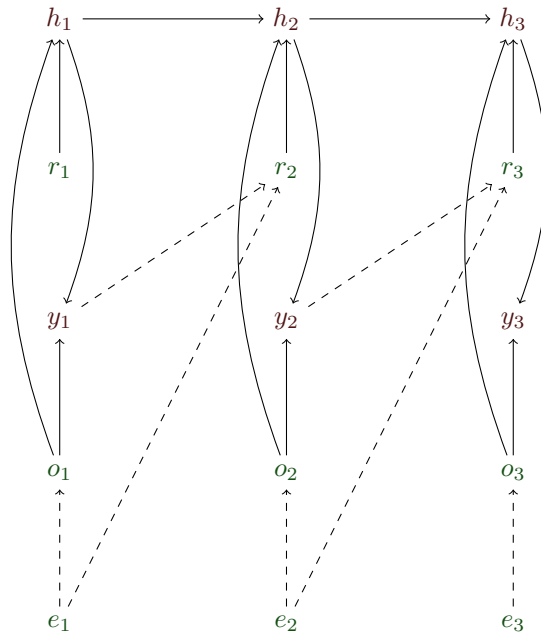


Figure 7: Influence diagram for hypothesis  $h_k$ , actions  $y_k$ , observations  $o_k$ , and rewards  $r_k$  for  $1 \leq k \leq 3$  with full lines corresponding to deterministic influences (via  $\pi$  and  $\mathcal{H}$ ) and dashed lines showing probabilistic influences (via  $\mu$  and  $\mathcal{E}$ ) in the nonsequential case.

## 1.8 Prior Knowledge

- *Implicit*: the setting of  $A$  (“hard bias”) and  $\mathcal{A}$  (“soft bias”)
- *Explicit*: the setting of  $a_1$  (“background knowledge”)

## 1.9 Hypothesis Representations

See Fig. 9.

## 1.10 Learning Scenarios

1. on-line concept learning
2. batch concept learning

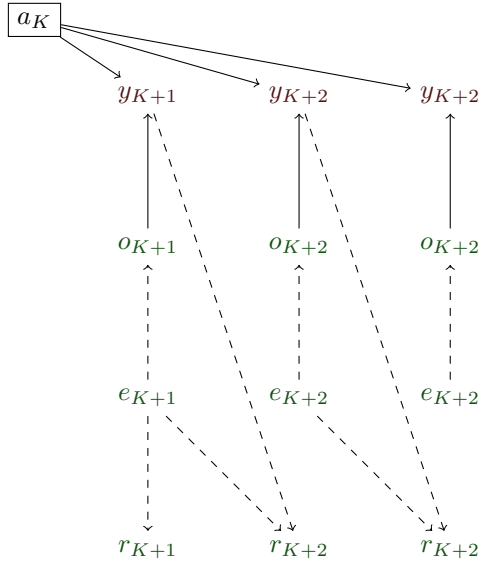


Figure 8: Influence diagram for actions  $y_k$ , observations  $o_k$ , states  $e_k$ , and rewards  $r_k$  in the action phase ( $k > K$ ) of batch learning with full lines indicating deterministic influences (via  $\pi$ ) and dashed lines showing probabilistic influences (via  $\mu$ ). The top row indicates the influence of the agent's last state in the training phase on the action phase. The dependence of  $r_{K+1}$  on  $e_K$  and  $y_K$  is not shown.

3. query-based and active learning (not covered here)
4. reinforcement learning
5. universal learning (not covered here)

## 2 On-line Concept Learning

We first motivate the on-line concept learning scenario with an example, in which the agent is an artificial scientist. The agent conducts repeated experiments with a living cell, which represents the environment. In each experiment, it observes two proteins of interest in the cell. More specifically, the agent detects whether the proteins are present in the cell at all, and it also determines whether they are in an active state (a special spatial conformation of a protein). The agent suspects that these proteins (both or only one of them) initiate apoptosis (cellular suicide). After each observation of the proteins, it tries to predict whether the cell will die or not. If the prediction is incorrect, the agent receives

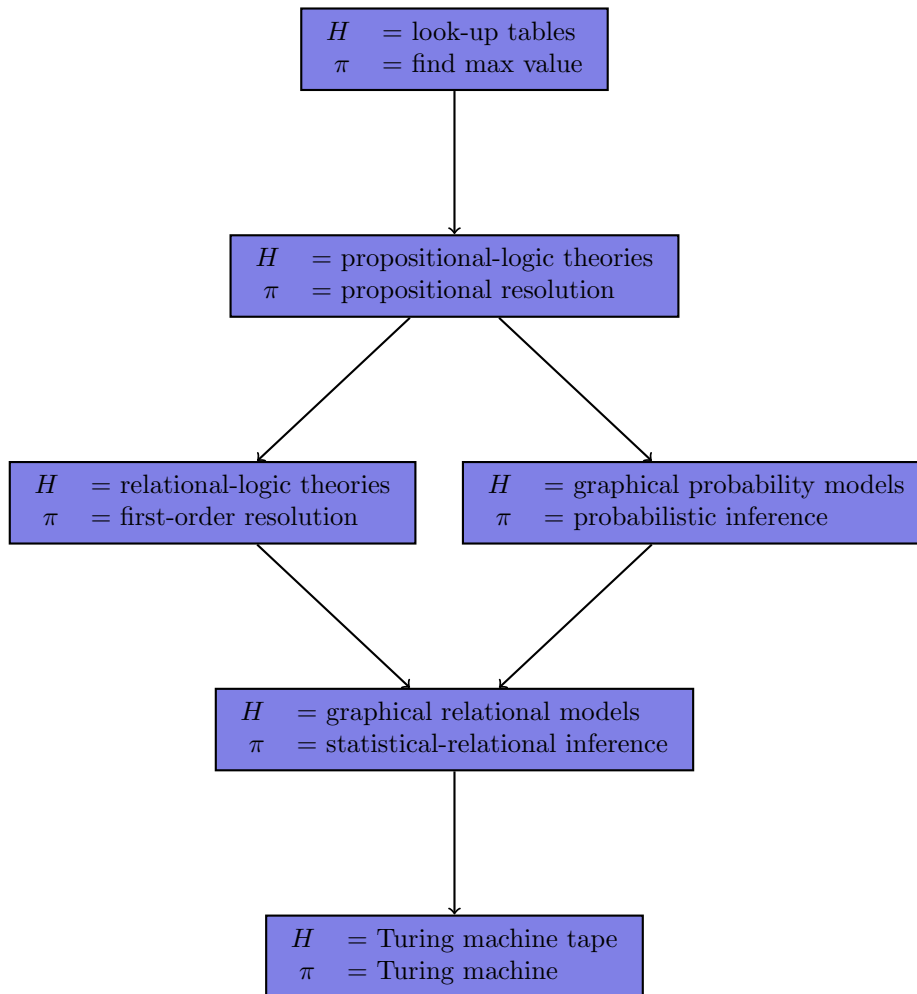


Figure 9: Hypothesis representations and their corresponding policy classes (interpreters) considered in this course. Arrow directions indicate increasing expressiveness.

a negative reward. This can be for example a cut-down on the agent's salary by the boss of the lab who is not happy with wrong biological predictions, in which case the boss would be a part of the environment. However, we will simply model such a reward with the number -1 for wrong predictions and with 0 for correct predictions.

Table 2 illustrates a history of such agent-environment interaction, in which the agent eventually learns that apoptosis is induced if and only if protein 1



| experiment number | apoptosis initiated | protein 1 present | protein 1 active | protein 2 present | protein 2 active | prediction | reward |
|-------------------|---------------------|-------------------|------------------|-------------------|------------------|------------|--------|
| $k$               | $e_k$               | $o_k^1$           | $o_k^2$          | $o_k^3$           | $o_k^4$          | $y_k$      | $r_k$  |
| 1                 | 0                   | 0                 | 0                | 0                 | 0                | 0          | 0      |
| 2                 | 0                   | 0                 | 0                | 1                 | 0                | 1          | 0      |
| 3                 | 0                   | 1                 | 0                | 0                 | 0                | 1          | -1     |
| 4                 | 1                   | 1                 | 0                | 0                 | 0                | 0          | -1     |
| 5                 | 0                   | 1                 | 0                | 1                 | 1                | 0          | -1     |
| 6                 | 1                   | 1                 | 1                | 1                 | 0                | 1          | 0      |
| 7                 | 1                   | 1                 | 1                | 0                 | 0                | 1          | 0      |
| 8                 | 0                   | 1                 | 0                | 1                 | 1                | 0          | 0      |
| (etc.)            |                     |                   |                  |                   |                  |            |        |

Table 1: An on-line concept learning experiment.

is present and it is in the active form. From time  $k = 5$  on, the agent makes correct predictions and is no longer punished with negative reward.

In the sequel, we will see how to model the illustrated scenario in our frameworks and we will see examples of agents able to learn as the agent-scientist has in the story above.

As Table 2 already indicated, the unknown variable guessed by the agent corresponds to the unknown state of the environment  $e$ . The variable is binary so we set

$$E = \{0, 1\} \tag{36}$$

We will accommodate (36) as a general assumption in the forthcoming text, unless we specify otherwise. This is because the binary setting is the simplest non-trivial one, to which richer state sets can usually be reduced.

The central assumption of concept learning is that the state  $e_k$  is fully *determined* by the observation  $o_k = (o_k^1, o_k^2, \dots)$ . In other words, an observation  $o$  generated by state 0 cannot be generated by state 1 and vice versa. In terms of the probability notation, this means that

$$\mu_O(o|0)\mu_O(o|1) = 0 \tag{37}$$

i.e., at most one of the probabilities is non-zero. Note that although the state is determined by the observation, the agent does not know *how* until it *learns* such knowledge.

The set of observations which can be generated from state 1 is called the *concept*  $C$  (of the environment)

$$C = \{o \in O \mid \mu_O(o|1) > 0\} \tag{38}$$

and the observations coming from this concept are *positive* observations (or, ‘examples’) of it, while the remaining observations, i.e. those in  $O \setminus C$  are termed *negative*. Since the environment states are partitioned into two *classes* (positive and negative), the agent’s guessing is an act of *classification* and the concept learning task is a special case of what is commonly termed *classification learning*.

As we have indicated already, the agent guesses the current state through its decision variable  $y \in Y$ . It is thus natural to set

$$Y = E = \{0, 1\} \quad (39)$$

Whenever the agent makes an incorrect guess  $y_k \neq e_k$ , it will receive a unit negative reward  $-1$  at the next time instant, so we instantiate (25), i.e.,  $\mu_R(r_{k+1}|e_k, y_k)$  to assign probability 1 to  $r_{k+1}$  such that

$$r_{k+1} = \begin{cases} 0 & \text{if } e_k = y_k \\ -1 & \text{otherwise} \end{cases} \quad (40)$$

Note that now the rewards are determined functionally rather than probabilistically, except for the first reward  $r_1$ , which is immaterial and is still sampled from the marginal  $\mu_R(r_1)$ .

In a more general setting, the unit punishment  $-1$  could be replaced by a value  $L(e_k, y_k)$  (called *loss*) which could be different for the two different cases of  $e_k \neq y_k$  possible in the present binary setting. Of course, in richer than binary settings, more different mistake kinds exist. We will not bother here with such generalizations.

The agent’s guesses are given by the policy (28) and at any time  $k$  they induce a subset of observations analogical to the unknown concept  $C$

$$C(a_k) = \{o \in O \mid \pi(a_k, o) = 1\} \quad (41)$$

$C(a_k)$  is the *agent’s concept* or the *hypothesized concept*. Note that we distinguish the concept  $C$  inherent to the environment for the agent’s concept  $C(a_k)$  at time  $k$  only by indicating the latter to be a function of the agent’s state  $a_k$ . To maximize rewards, the agent should evolve its state  $a_k$  so that it hypothesized concept co-incides with the target concept  $C$ , i.e. from some  $k$  on,

$$C(a_k) = C \quad (42)$$

If this is the case, the agent has *identified* the target concept.

A crucial question is *how* the agent’s state should be updated so that the guessing accuracy improves. A simple idea would be to let the agent remember a finite number of past observations and their true classes. When a new observation

comes, the agent would look up the most similar observation in this memory and guess the class associated with it. Here, similarity could be for example determined by the Hamming distance on the binary observation tuples.

We thus let the state act as a memory for a finite number  $m$  of observations and their true classes.

$$a_k = [(o_{k-m+1}, e_{k-m+1}), \dots, (o_{k-1}, e_{k-1}), (o_k, y_k)] \quad (43)$$

The true class  $e_k$  of observation  $o_k$  can be determined at time  $k + 1$  in the obvious way, given the guess  $y_k$  made for  $o_k$  and the reward  $r_{k+1}$  received for that guess. At time  $k$  the agent does not know  $r_{k+1}$  so for the newest observation it just stores its own guess  $y_k$  made according to the class of the most similar observation among  $o_{k-m+1}, \dots, o_{k-1}$ . In the state update step, the previous guess is replaced by the actual true state and the less recent item is discarded from  $a_k$ .

The similarity based approach just explained assumes that similar observations tend to have the same classes. Whether or not such an assumption is justified, the approach hardly merits to be called learning as it rests in plain memorization of observations. We would prefer an agent capable of *generalizing* the observation towards a *theory*, or *hypothesis* prescribing how observations determine environment states. Such a hypothesis can, for example, take the form of a set of logical rules, an equation, or a program in a programming language.

First consider that the agent's state is fully defined by the agent's current hypothesis, which we denote  $h_k$ , so

$$a_k = h_k \quad (44)$$

Then the decision policy (28) becomes

$$y_k = \pi(h_k, o_k) \quad (45)$$

and it is then natural to view  $\pi$  as an interpreter (a logical inference mechanism, an equation solver, a program interpreter, etc.) of  $h_k$ , while  $o_k$  plays the role of input data for  $h_k$ , according to which the decision is made. Similarly, (41) can be rewritten into

$$C(h_k) = \{o \in O \mid \pi(h_k, o) = 1\} \quad (46)$$

and the equality (47) is rephrased as

$$C(h_k) = C \quad (47)$$

As can be expected, the state update function (29) should change the last hypothesis  $h_{k-1}$  towards the current one  $h_k$  whenever a wrong guess  $y_{k-1}$  was

made (recall that due to (40) such a wrong guess is indicated to the agent by  $r_k = -1$ ). The change should lead to a correction of the hypothesis so that the same mistake does not happen again. To conduct such an update step at time  $k$ , we will need to know what the previous, wrongly classified observation  $o_{k-1}$  was. Since  $o_{k-1}$  is not an argument in (29), we need to memorize it within the agent state. This means we cannot get rid completely of a memory component of the agent state. So instead of (44) we rather consider the state to be a tuple consisting of the memorized previous observation and the current hypothesis

$$a_k = (o_k, h_k) \tag{48}$$

The update rule (29) then takes the more specific form

$$\mathcal{A}(a_{k-1}, x_k) = \mathcal{A}((o_{k-1}, h_{k-1}), (o_k, r_k)) = (o_k, h_k) \tag{49}$$

where  $h_k$  is determined from  $h_{k-1}, o_{k-1}$ , and  $r_k$  in a way depending on the particular learning strategy. We will visit some strategies in the coming sections.

While we needed the memorized observation  $o_k$  stored as part of the agent's state, this was just for the purpose of updating the hypothesis. To produce the decision  $y_k$ , the memorized observation is not needed so in the remainder of this chapter, we will still use the notation (45)-(47) including  $h_k$  rather than  $a_k = (o_k, h_k)$  as an argument.

For simplicity, we also will assume the observations to have binary components (as in the running example in Table 2), so, in general, they will  $n$ -tuples ( $n \in N$ ) from

$$O = \{0, 1\}^n \tag{50}$$

## 2.1 Generalizing Agent

Recall the example from Table 2 and observe that the components of each observation  $o_k$  (columns 2-4) correspond to logical propositions such as “Protein 1 is present” and “Protein 1 is active”, which we can denote  $p_1, \dots, p_4$  (one for each of the four observation columns), as customary in propositional logic. Given (50), the values  $o_k^1, \dots, o_k^4$  carry the logical (truth) values assigned to these respective propositions at time  $k$ .

An idea then suggests itself, that the agent's hypothesis  $h_k$  would be a propositional logic formula built with truth-valued variables  $p_1, \dots, p_4$ . Its truth value for the assignment  $o_k^1, \dots, o_k^4$  to the variables would determine the decision  $y_k$ . For example, the hypothesis that apoptosis initiates if and only if *protein 1 is present and it is active* would be written as

$$h_k = p_1 \wedge p_2 \tag{51}$$

If  $h_k$  is a logical conjunction, then the decision policy (45) takes the more specific form

$$y_k = \pi(h_k, o_k) = \begin{cases} 1 & \text{if } o_k \models h_k \\ 0 & \text{otherwise} \end{cases} \quad (52)$$

where  $o_k \models h_k$  means  $h_k$  is true given the truth-value assignments  $o_i$  to variables  $p_i$ ,  $1 \leq i \leq n$ . More precisely, we say that positive (negative, respectively) literal  $p_i$  ( $\neg p_j$ ) is *consistent* with observation  $o_k$  if  $o_k^i = 1$  ( $o_k^i = 0$ ). Then,  $o_k \models h_k$  holds if and only if all literals of conjunction  $h_k$  are consistent with  $o_k$ .

Let us design an agent that learns an unknown conjunction such as the above. The plan is to start with the most specific hypothesis (a conjunction of all literals, i.e. all propositional variables as well as their negations) and then successively delete all literals inconsistent with the received observations. So the initial hypothesis is gradually *generalized* towards the correct one.

In the present example, the agent has the initial hypothesis

$$h_1 = p_1 \wedge \neg p_1 \wedge p_2 \wedge \neg p_2 \wedge p_3 \wedge \neg p_3 \wedge p_4 \wedge \neg p_4 \quad (53)$$

This is the most specific hypothesis as it conjoins all possible conditions (literals). At the same time, this conjunction can of course never be true as it is self-contradictory. However, the agent's strategy is to successively remove from it all the literals that are inconsistent with the coming observations.

After the first percept has been received, i.e. for  $k > 1$ , the update rule (49) determines  $h_k$  according to

$$h_k = \begin{cases} h_{k-1} & \text{if } r_k = 0 \\ \text{delete}(h_{k-1}, o_{k-1}) & \text{otherwise} \end{cases} \quad (54)$$

where

$$\text{delete} \left( \bigwedge_{i \in I} p_i \bigwedge_{j \in J} \neg p_j, (o^1, o^2, \dots, o^n) \right) = \bigwedge_{\substack{i \in I \\ o^i = 1}} p_i \bigwedge_{\substack{i \in I \\ o^i = 0}} \neg p_i \quad (55)$$

So the *delete* function keeps exactly those literals from  $h_{k-1}$  which are consistent with  $o_{k-1}$ .

How do we know that through such an update rule the agent will indeed improve its guessing so that eventually it will only be receiving non-negative rewards?

First, we need to assume that there indeed exists a ‘correct’ conjunction  $h^*$ . It is correct in the sense that if  $h_k = h^*$ , then (47) holds. In other words,

$$\pi(h^*, o_k) = e_k, \forall o_k \in O \quad (56)$$

To resolve the question, we need a few lemmas.

**Lemma 2.1.**  $e_k = 1$  if and only if all literals of  $h^*$  are consistent with  $o_k$ .

The above lemma follows directly from (52) and (56).

**Lemma 2.2.** Whenever  $\text{delete}(h_{k-1}, o_{k-1})$  is called,  $e_{k-1} \neq y_{k-1}$ , and if  $e_{k-1} = 0$ , then all literals of  $h_{k-1}$  are consistent with  $o_{k-1}$ .

*Proof.* To see why Lemma 2.2 is true, note that according to (54),  $r_k \neq 0$  when  $\text{delete}$  is called. Due to (40), this means that  $e_{k-1} \neq y_{k-1}$ . So if  $e_{k-1} = 0$  then  $y_{k-1} = 1$ , but then due to (52),  $o_{k-1} \models h_{k-1}$  and so all literals of  $h_{k-1}$  are indeed consistent with  $o_{k-1}$ .  $\square$

**Lemma 2.3.**  $\text{delete}(h_{k-1}, o_{k-1})$  never removes a literal  $l \in h_{k-1}$  which is also in  $h^*$ .

*Proof.* Assume for contradiction that it does remove a literal  $l \in h^*$ . First assume  $e_{k-1} = 0$ . By Lemma 2.2, all literals of  $h_{k-1}$  are consistent with  $o_{k-1}$ . But because  $\text{delete}(h_{k-1}, o_{k-1})$  keeps all literals of  $h_{k-1}$  consistent with  $o_{k-1}$ , it does not delete  $l$ , which is a contradiction. Now consider  $e_{k-1} = 1$ . Then by Lemma 2.1 all literals of  $h^*$  including  $l$  must be consistent with  $o_{k-1}$ . Again, since  $\text{delete}$  keeps all consistent literals, it does not delete  $l$ , which is a contradiction.  $\square$

The starting hypothesis (53) of the designed agent is set to contain all possible literals, so  $h_1 \supseteq h^*$ , where the inclusion is with respect to the sets of literals in  $h_1$  and  $h^*$ . Furthermore, due to Lemma 2.3, we have

$$h_k \supseteq h^*, \forall k \in N \quad (57)$$

Given the above, the agent makes mistakes only on positive examples, and the mistakes are corrected by removing at least one inconsistent literal, as the following lemma formalizes.

**Lemma 2.4.** Assuming (53), whenever  $\text{delete}(h_{k-1}, o_{k-1})$  is called,  $e_{k-1} = 1$ , and the function deletes at least one literal from  $h_{k-1}$ .

*Proof.* Due to Lemma 2.2,  $e_{k-1} \neq y_{k-1}$ . If  $e_{k-1} = 0$  and  $y_{k-1} = 1$  then by the same lemma, all literals of  $h_{k-1}$  are consistent with  $o_{k-1}$ . According to Lemma 2.1,  $e_{k-1} = 0$  means that there is a literal in  $h^*$  inconsistent with  $o_{k-1}$ . But

due to (57), this inconsistent literal would also be contained in  $h_{k-1}$ , which is a contradiction. So we know that  $e_{k-1} = 1$  and  $y_{k-1} = 0$ . According to (52), this means that  $h_{k-1}$  contains a literal inconsistent with  $o_{k-1}$ . Since `delete`, by (55), keeps exactly all consistent literals, the inconsistent literal is removed.  $\square$

**Theorem 2.5.** *The generalizing agent makes at most  $2n$  mistakes, i.e. the cumulative reward is*

$$\sum_{k=1}^m r_k \geq -2n \quad (58)$$

for an arbitrary horizon  $m \in \mathbb{N}$ .

*Proof.* Since the first agent's conjunction has  $2n$  literals by (53) and upon each mistake, at least one literal is removed from the conjunction by Lemma 2.4, the maximum number of mistakes is  $2n$ .  $\square$

While the agent's strategy has been designed to learn conjunctions, it can be also made to learn disjunctions due to the equality

$$\neg(p_1 \vee p_2 \vee \dots \vee p_n) = \neg p_1 \wedge \neg p_2 \wedge \dots \wedge \neg p_n \quad (59)$$

So the only required change is that the agent replaces observations  $o_k$  with  $\bar{o}_k = (1 - o_k^1, 1 - o_k^2, \dots, 1 - o_k^n)$  and its actions  $y_k$  with  $1 - y_k$ .

Other logical classes can also be reduced to conjunction and disjunction learning. Consider e.g.  $s$ -CNF ( $s \in \mathbb{N}$ ). These are conjunctions of  $s$ -clauses. An  $s$ -clause is a disjunction of at most  $s$ -literals. There is a finite number of  $s$ -clauses so the agent can simply establish one new propositional variable for each possible  $s$ -clause a learn a conjunction with these new variables. This reduction would even be efficient if  $s$  is a small constant. Indeed, if  $n$  is the number of original variables, then the number of possible clauses is  $\binom{2n}{s}$ , i.e., the number of  $s$ -combinations of literals chosen from the set of  $n$  variables and their  $n$  negations. This number grows exponentially with  $s$  and polynomially with  $n$ . A similar reduction can be used to learn  $s$ -DNF.

## 2.2 The Subsumption Relation

It is instructive to view the generalization process as a path in the *subsumption lattice* of conjunctions shown for two propositional symbols in Fig. 10. A *lattice* is a partially ordered set where each two elements have their unique least upper bound and the greatest lower bound. The subsumption order is given by the subset relation

$$h_1 \subseteq h_2 \quad (60)$$

This means that conjunction  $h_1$  precedes conjunction  $h_2$  if the latter contains all literals of the former.

Recall from logic that a formula  $h_1$  *entails* another formula  $h_2$  if any model of  $h_1$  is also a model of  $h_2$ . We denote this as

$$h_1 \vdash h_2 \tag{61}$$

It is obvious that  $h_1 \subseteq h_2$  implies  $h_2 \vdash h_1$  if  $h_1$  and  $h_2$  are conjunctions. However, the inverse implication does not hold. For example (observe Fig. 10), we have both  $p_1 \wedge \neg p_1 \vdash p_2 \wedge \neg p_2$  and  $p_2 \wedge \neg p_2 \vdash p_1 \wedge \neg p_1$  simply because both of the formulas are non-satisfiable and thus neither has a model. However, they do not share any literal so the subset relation does not hold either way. Nevertheless, for satisfiable conjunctions (i.e., conjunctions other than ‘contradictions’)  $h_1, h_2, h_1 \subseteq h_2$  is equivalent to  $h_2 \vdash h_1$ .

While so far, we considered subsumption only conjunctions, the literal subset relation (60) is obviously defined as well for disjunctions, i.e. clauses. However, the relationship to logical entailment becomes inverted. More precisely, for two clauses  $h_1, h_2, h_1 \subseteq h_2$  implies  $h_1 \vdash h_2$ . Just like in the case of conjunctions, we cannot claim equivalence between the two latter relations. For example  $p_1 \vee \neg p_1 \vdash p_2 \vee \neg p_2$ . Again, the problem is with the atoms included both as a positive and a negative literals. While in conjunctions they produced contradictions, their presence in clauses make the latter tautologies, i.e. formulas true in any interpretation. But analogically to conjunctions,  $h_1 \subseteq h_2$  is equivalent to  $h_2 \vdash h_1$  if  $h_1, h_2$  are not tautologies.

Contradictory conjunctions and tautological clauses have one property in common. They contain a positive literal as well as the negation of the same literal. Clauses, which have this property, are called *self-resolving*.<sup>1</sup>

### 2.3 Separating agent

In Section 2.1 we have designed the generalization agent able to learn a conjunction while making only a finite number of mistakes. We have also seen that through efficient conversions, such an agent can also learn disjunction, *s*-DNF’s and *s*-CNF’s.

As an alternative example of a learning agent, we will demonstrate one using a different strategy to achieve the same goal. This time, the agent’s hypothesis  $h_k$  will be represented by non-logical means. In particular,  $h_k$  will define a hyperplane in the  $O = \{0, 1\}^n$  space (50) so  $C(h_k)$  (46) will include exactly those observations lying above the hyperplane.

<sup>1</sup>As the adjective *self-resolving* originates from the resolution principle, which is applied on clauses and not on conjunctions, it is usually not associated with conjunctions.



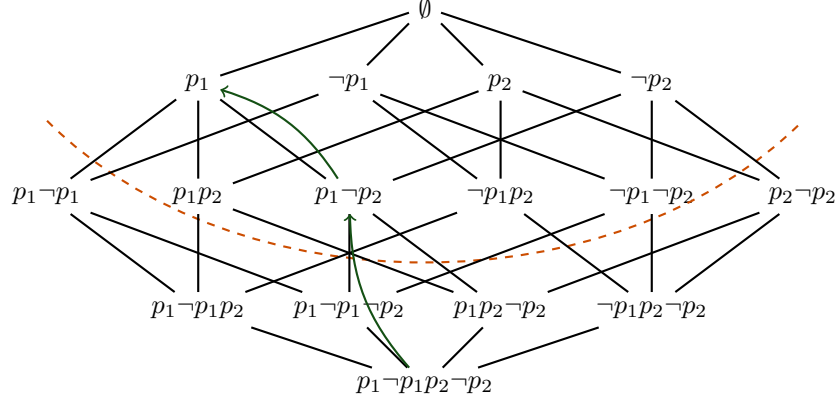


Figure 10: Subsumption lattice for conjunctions. The conjunction symbols  $\wedge$  are omitted for brevity. The curved arrows show how the agent generalizes its initial conjunction in two steps following the successive observations  $(1, 0)$  and  $(1, 1)$  carrying the respective truth values for  $p_1$  and  $p_2$ . All conjunctions below the dashed line are non-satisfiable.

Formally,  $h_k$  is an  $n$ -tuple of integer values bounded by some constant  $q \in N$ , i.e.  $h_k \in [0, 1, \dots, q]^n$ , so

$$h_k = [h_k^1, h_k^2, \dots, h_k^n] \quad (62)$$

The agent's decision policy (45) is given by a threshold function applied on a dot product

$$y_k = \pi(h_k, o_k) = \begin{cases} 1 & \text{if } h_k \cdot o_k > n/2 \\ 0 & \text{otherwise} \end{cases} \quad (63)$$

The initial hypothesis is

$$h_1 = (1, 1, \dots, 1) \quad (64)$$

And the update rule (49) is instantiated according to

$$h_k = \begin{cases} h_{k-1} & \text{if } r_k = 0 \\ \text{update}(2, h_{k-1}, o_{k-1}) & \text{if } h_{k-1} \cdot o_{k-1} \leq n/2 \\ \text{update}(0, h_{k-1}, o_{k-1}) & \text{if } h_{k-1} \cdot o_{k-1} > n/2 \end{cases} \quad (65)$$

wherein the function `update` is defined such that for  $h_k = \text{update}(\alpha, h_{k-1}, o_{k-1})$  and each  $i = 1, 2, \dots, n$ ,

$$h_k^i = \begin{cases} \alpha \cdot h_{k-1}^i & \text{if } o_{k-1}^i = 1 \\ h_{k-1}^i & \text{otherwise} \end{cases} \quad (66)$$

This agent, which can be considered an integer counterpart of the popular perceptron algorithm, learns a hyperplane. On the other hand, the generalization agent from the previous section was designed to learn logical formulas, namely conjunctions, disjunctions,  $s$ -DNF's, and  $s$ -CNF's. So how can we compare the two agents?

Assume that the target concept  $C$  corresponds to a disjunction  $c$  consisting of  $s$  literals made out of the variables  $p_1, \dots, p_n$ . That is to say,  $\mu_O(o|1) > 0$  if and only if  $o \models c$ . It is well known that disjunctions are linearly separable, so for a sufficiently large  $q$ , there is a hyperplane  $h^*$  such that (56) holds. This means, that the agent can identify a target disjunction through its hypothesis, although the latter is a hyperplane rather than a disjunction. The theorem below states that it does so with a finite number of mistakes.

**Theorem 2.6.** *The agent makes at most  $2+2s \lg n$  mistakes, i.e. the cumulative reward is*

$$\sum_{k=1}^m r_k \geq -2 - 2s \lg n \quad (67)$$

for any horizon  $m \in \mathbb{N}$ .

(proof omitted)

Just like the generalizing agent designed to learn conjunctions could easily be modified to learn disjunctions,  $s$ -CNF's, and  $s$ -DNF's, also the separating agent can be altered to learn conjunctions as well as the latter two classes by means of the same reduction principles. So the two agents can in principle learn the same concept classes. The difference is in the mistake bound. The latter agent performs better when the number of variables  $n$  is larger than the number of relevant variables  $s$ .

## 2.4 Hypothesis and Concept Classes

So far we have designed two exemplary learning agents, each with a different set of hypotheses  $h$  it could express. We shall call the set of all hypotheses an agent can express its *hypothesis class* denoted with the letter  $\mathcal{H}$  and we will assume  $\mathcal{H}$  to be finite. For the generalizing agent,  $\mathcal{H}$  consisted of all conjunctions made of at most  $n$  variables. For the separating agent,  $\mathcal{H}$  was the set of ( $q$ -bounded)  $n$ -tuples of integers.

Considering (46), each hypothesis class induces a set of concepts

$$\mathcal{C}(\mathcal{H}) = \{ C(h) \mid h \in \mathcal{H} \} \quad (68)$$

called a *concept class*.

The two agents exemplified so far used their update rules specifically designed for their respective hypothesis classes. Can we design a more general learning agent which could work with an arbitrary hypothesis class  $\mathcal{H}$ ?

Assume that  $\mathcal{H}$  is rich enough to contain a hypothesis  $h$  matching the unknown target concept  $C(h) = C$  (38). This assumption can be written as

$$C \in \mathcal{C}(\mathcal{H}) \tag{69}$$

Under such an assumption, since  $\mathcal{H}$  is finite, the agent can always try successively each element  $h \in \mathcal{H}$ , discarding it as soon as a mistake is made (negative reward received) using that hypothesis. In the worst case, the last hypothesis remaining will match the target concept. This means that the maximum number of mistakes made before identifying the target concept is

$$|\mathcal{H}| - 1 \tag{70}$$

This is a first indication of a dilemma we are going to face repeatedly in different forms. In particular, given that  $C$  is unknown, the agent should possess a large hypothesis space to maximize chances that (69) is satisfied. On the other hand, a large hypothesis space entails a loose bound on the number of mistakes made according to (70).

## 2.5 Version Space Agent

The general mistake bound (70) can be readily improved to  $\lg |C|$  using the *version space* strategy. Informally, its main idea is that on each observation, the agent discards all hypotheses from the hypothesis class which are inconsistent with the observation.

Before formalizing the principle, we will explain it by contrasting it to the generalization agent from Section 2.1. The latter agent worked with the hypothesis class  $\mathcal{H}$  of conjunctions made of at most  $n$  variables, i.e., at each time  $k$ , hypothesis  $h_k \in \mathcal{H}$ . The main idea of the version-space agent is to keep at each time  $k$  the subset of all hypotheses from  $\mathcal{H}$  which are consistent with all the observations received so far.

So instead of (71), the version-space agent's state is given as

$$a_k = (o_k, \mathcal{H}_k) \tag{71}$$

where  $o_k$  is again the memorized observation and  $\mathcal{H}_k$  is a set of hypotheses, also called the *version space*.

Decisions are determined by *voting* among all hypotheses in  $\mathcal{H}_k$ . Assuming that

$\mathcal{H}$  consists of logical formulas, this is formalized as

$$y_k = \pi(\mathcal{H}_k, o_k) = \begin{cases} 1 & \text{if } |\{ h_k \in \mathcal{H}_k \mid o_k \models h_k \}| > |\mathcal{H}_k|/2 \\ 0 & \text{otherwise} \end{cases} \quad (72)$$

The agent starts with the full hypothesis space

$$\mathcal{H}_1 = \mathcal{H} \quad (73)$$

The hypothesis update step takes a form slightly different from (49), in particular

$$\mathcal{A}(a_{k-1}, x_k) = \mathcal{A}((o_{k-1}, \mathcal{H}_{k-1}), (o_k, r_k)) = (o_k, \mathcal{H}_k) \quad (74)$$

where  $\mathcal{H}_k$  is obtained by deleting all hypotheses inconsistent with  $o_{k-1}$ , i.e.

$$\mathcal{H}_k = \begin{cases} \{ h \in \mathcal{H}_{k-1} \mid o_{k-1} \models h \} & \text{if } e_{k-1} = 1 \\ \{ h \in \mathcal{H}_{k-1} \mid o_{k-1} \not\models h \} & \text{if } e_{k-1} = 0 \end{cases} \quad (75)$$

where  $e_{k-1}$  is determined as  $e_{k-1} = |y_{k-1} - r_{k-1}|$  (check that this is true) and  $y_{k-1} = \pi(\mathcal{H}_{k-1}, o_{k-1})$ . To see that the latter expression can be evaluated at the hypothesis update step, note that both  $\mathcal{H}_{k-1}$  and  $o_{k-1}$  are indeed available at that step according to (75).

In (72) and (75) we used the relation  $\models$  assuming that  $h_k$  is a logical formula and  $o_k$  provides the truth-values for its atoms (variables). If this was not the case, we could assume more generally that hypotheses provide a mapping  $h_k : O \rightarrow \{0, 1\}$ . Then in the two formulas,  $o_k \models h_k$  ( $o_k \not\models h_k$ , respectively) would be replaced by  $h_k(o_k) = 1$  ( $h_k(o_k) = 0$ ).

Assume that  $\mathcal{H}$  is rich enough so that (69) is true. Then the following theorem holds.

**Theorem 2.7.** *The agent makes at most  $\lg |\mathcal{H}|$  mistakes, i.e. the cumulative reward is*

$$\sum_{k=1}^m r_k \geq -\lg |\mathcal{H}| \quad (76)$$

for any horizon  $m \in \mathbb{N}$ .

*Proof.* To see why the theorem holds note that the agent decides by the majority of hypotheses from  $\mathcal{H}_k$ . So if a mistake is made, at least half of the hypotheses in  $\mathcal{H}_k$  are deleted. In the worst case, the last remaining hypothesis is correct.  $\square$

Once again, a dilemma is observed in that  $\mathcal{H}$  should be large enough so that (69) is satisfied. However, the size  $|\mathcal{H}|$  also increases the mistake bound (76). The latter mistake bound is logarithmic, which is certainly a significant improvement over (70) good but the computational demands for storing  $\mathcal{H}_k$  (containing a potentially large number of hypotheses) can be prohibitive.

## 2.6 The Mistake Bound Learning Model

The linear mistake bounds we obtained for the generalizing and separating agents indicate that these agents are indeed able to learn well the conjunctive and disjunctive concepts but also other kinds of concepts (namely,  $s$ -DNF and  $s$ -CNF) that can be reduced to the latter. We will now generalize the notion of ‘good on-line learning.’ We say that an agent *learns class  $\mathcal{H}$  on-line* if it makes at most  $\text{poly}(n)$  of mistakes in the on-line scenario provided that (69) holds. Here,  $\text{poly}$  is a polynomial and  $n$  is the size of observations. With our setting (50), the size of observations is the number  $n$  of binary values making up the observations.

By Theorems 2.5 and 2.6, the generalizing and separating agents learn  $\mathcal{H}$  on-line when  $\mathcal{H}$  are conjunctions, disjunctions,  $s$ -DNF or  $s$ -CNF (for a constant  $s$ ). By Theorem 2.7, for an arbitrary class  $\mathcal{H}$ , the version-space algorithm has a mistake bound  $\lg |\mathcal{H}|$  as long as (69) holds true. If furthermore  $|\mathcal{H}|$  is at most exponential in  $n$ , the agent necessarily learns  $\mathcal{H}$  on-line, because the mistake bound  $\lg |\mathcal{H}|$  is then polynomial. The condition ‘at most exponential’ above seems rather permissive. But note that  $|\mathcal{H}|$  may easily be super-exponential. The extreme example of such a situation is the class  $\mathcal{H}$  such that (69) is true for any possible concept  $C$ ,  $C \subseteq O$ . Since  $|O| = |\{0,1\}^n| = 2^n$ , we have  $|\mathcal{H}| \geq 2^{|O|} = 2^{2^n}$ , so  $|\mathcal{H}|$  is super-exponential.

Furthermore, we refine the definition into a stricter form. An agent that learns hypothesis class  $\mathcal{H}$  on-line is said to learn it *efficiently* if it spends at most polynomial time (in  $n$ ) between the receipt of a percept and the generation of the next action.

It is quite easy to verify that both the generalizing and separating agents learn their hypothesis classes efficiently. However, the version-space agent is evidently not efficient if  $|\mathcal{H}|$  is super-polynomial in  $n$  as it has to visit each element of  $\mathcal{H}$  during the update (75). Say  $\mathcal{H}$  are conjunctions. There are  $2n$  literals, each of which can be present or absent in a conjunction so there are  $|\mathcal{H}| = 2^{2n}$  conjunctions and the version-space agent does not learn conjunctions efficiently. We could reduce the size of  $\mathcal{H}$  realizing that conjunctions containing the same atom both in a positive literal and a negative literal are unsatisfiable and can be discarded. Then we would derive  $|\mathcal{H}|$  as follows. We have  $n$  atoms, each of which may be positive, negative, or absent in a conjunction. Thus there are  $3^n$  conjunctions in total, which grows slower than  $2^{2n}$  but still exponentially.

What about a *lower bound* on mistakes? The latter can be established using the notion called *VC-dimension* of a hypothesis class. We say that a set of observations  $O' \subseteq O$  is *shattered* by hypothesis class  $\mathcal{H}$  if

$$\{ O' \cap C(h) \mid h \in \mathcal{H} \} = 2^{O'} \quad (77)$$

which means that the set of observations can be partitioned in all possible ways into two classes by the hypotheses from  $\mathcal{H}$ . The *Vapnik-Chervonenkis Dimension* (or VC-dimension) of  $\mathcal{H}$ , written  $\text{VC}(\mathcal{H})$ , is the cardinality of the largest set  $O' \subseteq O$  that is shattered by  $\mathcal{H}$ .

**Theorem 2.8.** *Assume an agent deciding by (45) where  $h_k \in \mathcal{H}$ . No upper bound on the number of mistakes made by the agent is smaller than  $\text{VC}(\mathcal{H})$ .*

*Proof.* For any sequence of agent's decisions  $y_1, y_2, \dots, y_{\text{VC}(\mathcal{H})}$  there exists a  $h \in \mathcal{H}$  according to which all these decisions are wrong.  $\square$

The theorem does not directly apply to the version-space agent, which decides according to (72) rather than (45). This significance of  $\text{VC}(\mathcal{H})$  for this agent is that it will not identify the target hypothesis (i.e. reduce  $|\mathcal{H}_k|$  to 1) in fewer than  $\text{VC}(\mathcal{H})$  steps. Indeed,  $\mathcal{H}_{\text{VC}(\mathcal{H})}$  will contain two hypotheses, one deciding 1 for  $o_{\text{VC}(\mathcal{H})}$  and another deciding otherwise.

### 3 Batch Concept Learning

We will maintain the assumption of concept learning (37), that is, no observation can be generated from two different states of the environment. For simplicity, we will continue working in the binary setting (39), i.e. the environment has only two different states and the agent can make only two different actions. Also, the observations are still assumed to be binary (50). Finally, rewards are still determined according to (40), i.e.,  $-1$  is a punishment for a wrong guess.

However, instead of the on-line setting, we will now study concept learning in the batch framework we defined in Section 1.7. Since observations are i.i.d. here, we can express the probability distribution of reward  $r_{k+1}$  by (34) as  $\mu_r(r_{k+1}|a_k)$ , only depending on the current agent's state. Considering that the decision for an observation  $o_k$  is made only using the hypothesis component  $h_k$  of the state  $a_k$  (45), this is the same as writing  $\mu_r(r_{k+1}|h_k)$ . As follows from (40),  $\mu_r(-1|h_k)$  is the probability that the decision  $y_k = \pi(h_k, o_k)$  is wrong, i.e.,  $y_k \neq e_k$ . We will call this probability the *error* of  $h_k$

$$\text{err}(h_k) = \mu_r(-1|h_k) \tag{78}$$

Remind that the agent's goal in the batch setting is to reach at the end of the training phase (i.e., at time  $K$ ), a state  $a_K$  which maximizes the expected reward (35). Again, since only the hypothesis component of the agent's state determines the decision, the expected reward is  $\sum_{r_k \in R} \mu_R(r_k|h_K)r_k$ . This expectation is

equal to  $-\text{err}(h_K)$  as can be easily seen realizing that rewards can be only 0 or -1.

So the goal of batch learning can be equivalently stated as arriving at a hypothesis  $h_K$  minimizing the error. A natural question of interest is how the algorithms we designed for on-line concept learning in the sequential scenario would perform in terms of the error. Evidently, the bounds on the number of mistakes we established in Theorems 2.5, 2.6, and 2.7 do not translate to any bound on  $\text{err}(h_K)$  as there is no guarantee that the mistakes will happen in the learning phase ( $k \leq K$ ) where the agent still can fix its hypothesis.

But unlike in the on-line learning case, the batch case inherits the non-sequential assumptions (30) and (31), meaning that states and observations are sampled i.i.d. according to distributions that do not change with  $k$ . They prevent the environment from ‘adversarial’ behavior, for example, one where the training phase would only contain ‘easy’ examples and the ‘hard’ ones would be kept for the testing phase. As we will see, in this scenario we can indeed bound  $\text{err}(h_K)$  for particular learning agents, although we will be able to do it only with certain probability smaller than 1.

### 3.1 Batch Learning with the Generalizing Agent

Consider the generalizing agent as described in Section 2.1 working in the learning phase ( $k \leq K$ ) of the batch scenario just as it worked in the on-line scenario.

Denote by  $\Pr(l)$  the probability that literal  $l \in h_{k-1}$  is deleted by (55) while updating to  $h_k$ . Since the agent makes mistakes only on positive examples, this means that  $o_{k-1}$  was positive and  $l$  was inconsistent with it. Observations are i.i.d. in the non-sequential batch setting so  $\Pr(l)$  does not depend on  $k$ .

According to (54), deletion takes place if and only if  $r_{k-1} = -1$ . So considering (78),  $\text{err}(h_k)$  is equal to the probability that *some* of the literals in  $h_{k-1}$  get deleted, which in turn can be bounded by the sum

$$\text{err}(h_k) \leq \sum_{l \in h_k} \Pr(l) \tag{79}$$

We have no more than  $2n$  literals in  $h_k$  so if  $\Pr(l) \leq \epsilon/2n$  ( $\epsilon \in \mathfrak{R}$ ) for each of them then  $\text{err}(h_k) \leq \epsilon$ . Let us fix an arbitrary  $\epsilon$  and let us call a literal *bad* if  $\Pr(l) > \epsilon/2n$ .

Consider a bad literal  $l$  from  $h_1$  (53). The probability that it will still be in  $h_2$  (i.e., that it will ‘survive’ the first observation) is  $1 - \Pr(l)$ . Appealing again to

the i.i.d. character of observations, the probability of  $l \in h_{k+1}$  is then

$$(1 - \Pr(l))^k < \left(1 - \frac{\epsilon}{2n}\right)^k \quad (80)$$

There are at most  $2n$  bad literals so the probability that *some* of them is in  $h_{k+1}$  is at most

$$2n \left(1 - \frac{\epsilon}{2n}\right)^k \quad (81)$$

To work with this upper bound easily, we make use of the inequality  $1 - x \leq e^{-x}$  which holds for  $x \in [0; 1]$ , to obtain

$$2n \left(1 - \frac{\epsilon}{2n}\right)^k \leq 2ne^{-k\frac{\epsilon}{2n}} \quad (82)$$

We now summarize the above inferences into a theorem.

**Theorem 3.1.** *Hypothesis  $h_{k+1}$  of the generalizing agent in the learning phase ( $k < K$ ) has  $\text{err}(h_{k+1}) \leq \epsilon$  with probability at least  $1 - 2ne^{-k\frac{\epsilon}{2n}}$*

So at the end of learning,  $\text{err}(h_K) < \epsilon$  with probability at least  $1 - 2ne^{-m\frac{\epsilon}{2n}}$  where  $m = K - 1$ .

## 3.2 Batch Learning with General On-line Agents

We define a *standard* on-line agent as one that uses a single hypothesis to determine its decisions (45), and changes the hypothesis if and only if a mistake has been made by the previous hypothesis. This includes the generalizing and separating agents as follows from the update rules (54) and (65). On the other hand, the version-space agent is not standard for obvious reasons.

The next lemma will enable us to accommodate any standard on-line learning agent for the batch learning scenario with a probabilistic bound on the error of the learned hypothesis.

**Lemma 3.2.** *If a standard on-line agent retains a hypothesis  $h_k$  for  $q$  steps ( $h_k = h_{k+1} = \dots = h_{k+q}$ ), then  $\text{err}(h_k) \leq \epsilon$  with probability at least  $1 - e^{-q\epsilon}$ .*

*Proof.* The probability that the standard agent retains a bad hypothesis (i.e., one for  $\text{err}(h_k) > \epsilon$ ) on receiving an observation is the probability  $1 - \text{err}(h_k)$  that the bad hypothesis produces a correct decision for that observation. Since  $\text{err}(h_k) > \epsilon$ , the probability is at most  $1 - \epsilon$ . The probability of keeping the hypothesis over  $q$  i.i.d. observations is thus at most  $(1 - \epsilon)^q$ , and we already know that  $(1 - \epsilon)^q \leq e^{-q\epsilon}$ . Otherwise, i.e. with probability at least  $1 - e^{-q\epsilon}$ , the hypothesis was not bad, i.e.  $\text{err}(h_k) \leq \epsilon$ .  $\square$



So the rule is: wait until  $h_k = h_{k+1} = \dots h_{k+q}$  happens and then keep  $h_k$  with the probabilistic error bound. The question is how to guarantee that the event indeed happens within the learning phase, i.e.  $k + q \leq K$ . If we have a mistake bound  $M$  for the agent, we know that the standard agent makes at most  $M$  hypothesis changes. In this case we set the learning phase long enough, in particular  $K = Mq$ , to guarantee that one of the hypotheses in the learning phase survives at least  $q$  observations.

### 3.3 Consistent Agent

Here we design a general agent working with an arbitrary hypothesis space. This is analogical to the version space agent we studied in the on-line setting.

We first adapt the version space agent from on-line learning to batch learning. In the learning phase, the agent works just as in the on-line setting. When the phase ends, i.e.  $k = K$ , the agent updates its state for the last time according to (75), and then selects an *arbitrary* hypothesis  $h_K$  from  $\mathcal{H}_K$ .

Once again, let us call a hypothesis  $h \in \mathcal{H}$  *bad* if  $\text{err}(h) > \epsilon$ . The probability that a particular bad hypothesis  $h \in \mathcal{H}$  is in  $h_K$  is at most  $(1 - \epsilon)^m \leq e^{-\epsilon m}$  where<sup>2</sup>  $m = K - 1$ . The probability that *some* bad hypothesis from  $\mathcal{H}$  (73) is in  $h_K$  is at most  $|\mathcal{H}|e^{-\epsilon m}$ . So the probability that no bad hypothesis survives and thus  $\text{err}(h_K) \leq \epsilon$  whichever  $h_K$  the agent has picked from  $\mathcal{H}_K$ , is at least  $1 - |\mathcal{H}|e^{-\epsilon m}$ .

Maintaining the sets  $\mathcal{H}_1, \mathcal{H}_2, \dots$  along the course of learning can of course be prohibitive as the latter can be very large sets. A behavior equivalent to the described version-space strategy can however be obtained without maintaining the version spaces.

In particular, the *consistent agent* collects all observations seen during the training phase along with their true classes in its memory without updating its hypothesis. Only at the end of the training phase, it picks from its hypothesis class a hypothesis consistent with all the collected observations.

Formally, the agent has state

$$a_k = ((o_1, e_1, \dots, o_{k-1}, e_{k-1}), h_k) \tag{83}$$

for time  $k = 1, 2, \dots, K - 1$ . Here,  $h_k$  is just any ‘dummy’ hypothesis that makes  $y_k = \pi(h_k, o_k) = 0$  for any observation  $o_k \in O$  (e.g., an unsatisfiable

---

<sup>2</sup>Remind that  $o_K$  is the first testing observation (its true class  $e_k$  cannot be determined at time  $K$ ). So  $m = K - 1$  is the number of observations with known class (label) received in the training phase. In other words,  $m$  is the number of labeled training examples received by the agent.

formula). With  $y_k = 0$  the agent can indeed determine the true classes in (83) as  $e_k = -r_{k+1}$ . Only at time  $K$ , the agent finds any hypothesis  $h_K \in \mathcal{H}$  consistent with the collected set, i.e.  $\pi(h_K, o_k) = e_k$  for all  $k < K$ . Analogically, to the reasoning above, we have that

**Lemma 3.3.** *The probability that the consistent agent’s hypothesis  $h_K$  has error  $\text{err}(h_K) \leq \epsilon$  is at least  $1 - |\mathcal{H}|e^{-\epsilon m}$  where  $m = K - 1$ .*

The first component of the agent’s state at time  $K$ , i.e.

$$o_1, e_1, \dots, o_m, e_m \tag{84}$$

is called the *training set*. Unlike the version space agent, the consistent agent does not store large hypothesis spaces and so it obviously consumes less memory. On the other hand, finding a hypothesis  $h_K$  consistent with the training set may be computationally hard, and potentially intractable.

### 3.4 The PAC Learning Model

Agent *probably approximately learns class  $\mathcal{H}$  (in the batch setting)* if at the end of the training phase it produces  $h_K$  such that  $\text{err}(h_K) \leq \epsilon$  with probability at least  $1 - \delta$ , and  $K - 1 = m \leq \text{poly}(n, 1/\delta, 1/\epsilon)$ , where *poly* is a polynomial.

“probably approximately learns” = “PAC-learns” (C for correctly)

It PAC-learns the class *efficiently* if it spends at most polynomial (in the same variables) time between the receipt of a percept and the generation of the next action in the training phase.

**Theorem 3.4.** *The generalizing agent efficiently PAC-learns conjunctions.*

*Proof.* Efficiency is obvious: at most  $2n$  unit steps (going over literals) for each of  $m = K - 1$  observations. From Theorem 3.1, the probability that  $\text{err}(h_K) > \epsilon$  is at most  $2ne^{-m\frac{\epsilon}{2n}}$ . It remains to determine how many observations  $m$  are needed to make the probability smaller than a given  $\delta$  and see if the result is polynomial.

$$\begin{aligned} \delta &> 2ne^{-m\frac{\epsilon}{2n}} \\ \frac{\delta}{2n} &> e^{-m\frac{\epsilon}{2n}} \\ \ln \frac{\delta}{2n} &> -m\frac{\epsilon}{2n} \\ \frac{2n}{\epsilon} \ln \frac{2n}{\delta} &\leq m \end{aligned}$$

So the required  $m = K - 1$  is indeed polynomial in  $n$ ,  $1/\epsilon$  and  $1/\delta$ .  $\square$

**Theorem 3.5.** *Any standard agent learning (efficiently) a concept class  $C$  on-line, has a counterpart which (efficiently) PAC-learns  $C$ .*

*Proof.* The agent makes at most  $M < \text{poly}(n)$  updates, i.e. max number of mistakes.

Its batch counterpart works as follows.

$$\text{Set } q = \frac{1}{\epsilon} \ln\left(\frac{1}{1-\delta}\right)$$

If before  $M$  updates have been made, each hypothesis survived for less than  $q$  steps, then the last one (which makes no mistakes) is found in at most  $Mq$  steps, and is kept as  $h_K$ . Both  $M$  and  $q$  are polynomial.

If some of them survived for at least  $q$  steps, than according to lemma (3.2), its error is less than  $\epsilon$  with probability at least  $1 - e^{-q\epsilon} = \delta$ . This hypothesis found with less than  $Mq$  (poly) steps, will be kept as  $h_K$ .  $\square$

So a negative batch (PAC) result also means a negative on-line result.

**Theorem 3.6.** *If (69) holds and  $|\mathcal{H}|$  is at most exponential in  $n$  then the consistent agent PAC-learns  $\mathcal{H}$ .*

*Proof.* By Lemma (3.3), probability  $\delta$  that  $\text{err}(h) > \epsilon$  is at most  $|\mathcal{H}|e^{-\epsilon m}$ .

$$\delta \leq |\mathcal{H}|e^{-\epsilon m}$$

$$\frac{\delta}{|\mathcal{H}|} \leq e^{-\epsilon m}$$

$$\frac{1}{\epsilon} \ln \frac{|\mathcal{H}|}{\delta} > m$$

Since  $|\mathcal{H}|$  is at most exponential in  $n$ ,  $\ln |\mathcal{H}|$  is at most polynomial in it, so  $m < \text{poly}(1/\epsilon, 1/\delta, n)$ .  $\square$

Also,  $s$ -CNF and  $s$ -DNF learnable by poly reduction to conjunctions.