

# Computational Intelligence



LMU Munich  
winter term 2024/2025

Thomas Gabor  
Claudia Linnhoff-Popien

# reading exercise

## 2024-11-28

### TECHNICAL ARTICLES

THE PROLOG PREDICATOR  
Dave McInnis  
This University

In 1961, Carl Hewitt introduced **PLANNER**, a procedural deductive system. [Hewitt 72] It incorporated the very innovative concept for the procedural interpretation of that time, the idea of *goals* or *assertions*, as in the induction, deduction, and resolution of logic, and was implemented by McCarthy, Minsky, and Charniak [McCarthy 71] and used by Minsky in his excellent book on natural language. [McCarthy 71] A very similar language, **GAL**, was implemented in the next year. [Hewitt 72] These languages, known by the name of *goal languages*, were widely used as providing facilities for the next generation of AI systems, the way LISP had provided facilities for the previous one.

By 1971, however, people were being dissatisfied with these languages because known the 711 some attempts were made to implement improved systems without the drawbacks [McCarthy 71], but they didn't really catch on. By 1975, no one was using the "AI languages" in 1971, so one was forced to use it to implement tools that required it. This was the language of choice, and people started to use it to implement tools that required it. This was the language of choice, and people started to use it to implement tools that required it.

At about the time AI languages were dying [Hewitt 72] and Robert Hewitt (now John AI) discovered the procedural interpretation of **PROLOG** (for "Procedural Logic") that allowed them to use it to implement tools that required it. This was the language of choice, and people started to use it to implement tools that required it.

This has not happened. **PROLOG** has attracted devoted as much attention as LISP, and it is not surprising that it is growing. In a recent issue of **PROLOG** with some success of this language, and powerful language, and it is growing. In a recent issue of **PROLOG** with some success of this language, and powerful language, and it is growing. In a recent issue of **PROLOG** with some success of this language, and powerful language, and it is growing.

In LISP, you distinguish between a program and a single function. In **PROLOG**, instead of functions you have *relations*. A relation is an ordered set of elements. A clause is of the form

pattern :- body

meaning, "to match pattern, do body." In the case of *literal*, of the form *literal*, although *literal* is allowed, it will follow.

For example, we can say

```
quadrant(A,B,C,Walroots) :-  
  disjunct(A,B,C,D) :-  
    mult(A,B,C,D), mult(A,C,D),  
    mult(A,D,C,D), mult(A,D,C,D).
```

quadrant(A,B,C,D) :- no.

```
quadrant(A,B,C,D) :-  
  disjunct(A,B,C,D) :-  
    mult(A,B,C,D), mult(A,C,D),  
    mult(A,D,C,D), mult(A,D,C,D).
```

```
quadrant(A,B,C,D) :-  
  disjunct(A,B,C,D) :-  
    mult(A,B,C,D), mult(A,C,D),  
    mult(A,D,C,D), mult(A,D,C,D).
```

This program finds the end points of a line segment. It does this by using the traditional language. The first thing it does is to find the end points of the line segment. It does this by using the traditional language. The first thing it does is to find the end points of the line segment.

I just quoted some of the last lines of **PROLOG** before. Now, I show you some of the first lines of **PROLOG**. I show you some of the first lines of **PROLOG**. I show you some of the first lines of **PROLOG**. I show you some of the first lines of **PROLOG**.

The old relation is a primitive relation that is not a function. It is a relation that is not a function. It is a relation that is not a function. It is a relation that is not a function.

### Natural Language Processing With Prolog in the IBM Watson System

Adam Lally  
IBM Thomas J. Watson Research Center  
Paul Fodor  
Stony Brook University  
24 May 1971

On February 14-16, 1971, the IBM Watson question answering system won the Jennings and Brad Hunt. To compete successfully at Jennings, Watson had to answer complex natural language questions over an extremely broad domain of knowledge. Moreover, it had to compute an accurate confidence in its answers and to complete its processing in a very short amount of time.

The Question Answering (QA) problem requires a machine to go beyond just matching key words in documents, which is what a web search engine does, and correctly interpret the question to figure out what is being asked. The QA system also needs to find the precise answer without requiring the aid of a human to read through the returned documents.

To address these challenges, the research team at IBM developed a software architecture called DeepQA, on which Watson is implemented. The DeepQA architecture assumes and proves multiple interpretations of the question, generates many plausible answers or hypotheses, collects evidence for these hypotheses, and evaluates the evidence to determine if it supports or refutes these hypotheses [2]. Watson contains hundreds of different algorithms that evaluate evidence along different dimensions.

Watson utilizes Natural Language Processing (NLP) technology to interpret the question and extract key elements such as the answer type and relationships between entities. Also, NLP was used to analyze (prior to the competition) the vast amounts of unstructured text (encyclopedias, dictionaries, news articles, etc.) that may provide evidence whether the relationships between entities in the question match those in the source.

Watson's NLP begins by applying a parser [5] that corrects each text sentence into a more structured form, a tree that shows both surface structure and deep, logical structure. For example, in the following example Jennings' question:

POETS & POETRY: He was a hard clerk in the Yukon before he published "Songs of a Seafarer" in 1907

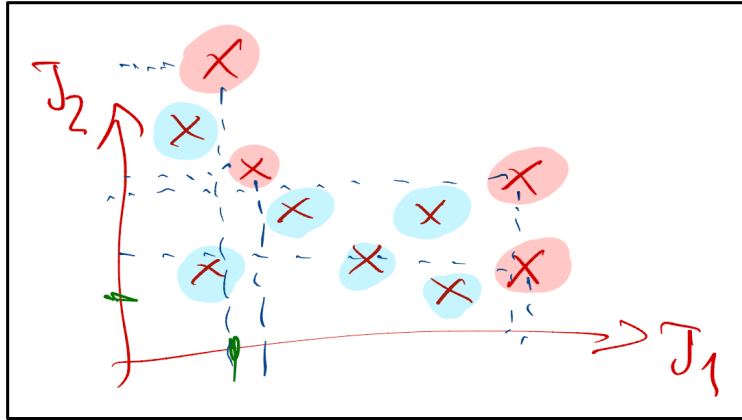
“early exam”

2025-02-06

18:30h

on main campus

# Multi-Objective Optimization



**Definition 6** (multi-objective optimization). Let  $\mathcal{E} = (\mathcal{X}, \mathcal{T}, \tau, E, \langle X_u \rangle_{0 \leq u \leq t})$  be an optimization process.  $\mathcal{E}$  is a multi-objective optimization process iff the target space  $\mathcal{T}$  has the form  $\mathcal{T} = \mathcal{T}_0 \times \cdots \times \mathcal{T}_{N-1}$  for some  $N \in \mathbb{N}$  with  $\leq_i$  being a total order on  $\mathcal{T}_i$  for any  $i \in [0; N-1] \subset \mathbb{N}$ . Unless stated otherwise, we assume that no single total order on  $\mathcal{T}$  is available. However, we can construct a partial order  $\preceq$  so that

$$(x_0, \dots, x_{N-1}) \preceq (x'_0, \dots, x'_{N-1}) \iff \forall i \in [0; N-1] \subset \mathbb{N} : x_i \leq x'_i,$$

which is sufficient to adapt many standard optimization algorithms.

**Definition 6** (multi-objective optimization). Let  $\mathcal{E} = (\mathcal{X}, \mathcal{T}, \tau, E, \langle X_u \rangle_{0 \leq u \leq t})$  be an optimization process.  $\mathcal{E}$  is a multi-objective optimization process iff the target space  $\mathcal{T}$  has the form  $\mathcal{T} = \mathcal{T}_0 \times \cdots \times \mathcal{T}_{N-1}$  for some  $N \in \mathbb{N}$  with  $\leq_i$  being a total order on  $\mathcal{T}_i$  for any  $i \in [0; N-1] \subset \mathbb{N}$ . Unless stated otherwise, we assume that no single total order on  $\mathcal{T}$  is available. However, we can construct a partial order  $\preceq$  so that

$$(x_0, \dots, x_{N-1}) \preceq (x'_0, \dots, x'_{N-1}) \iff \forall i \in [0; N-1] \subset \mathbb{N} : x_i \leq x'_i,$$

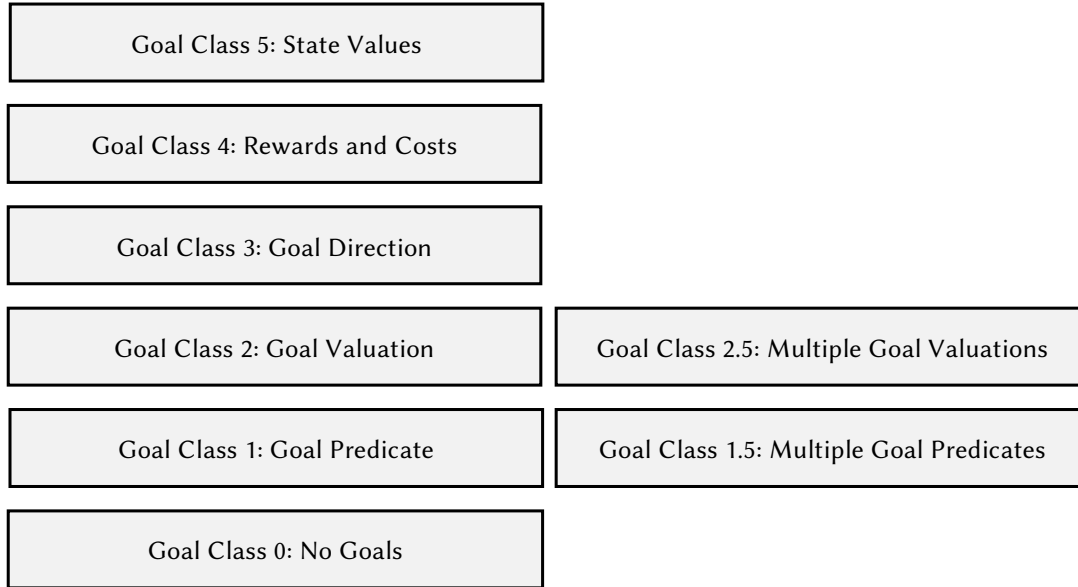
which is sufficient to adapt many standard optimization algorithms.

**Definition 7** (Pareto front for optimization). Let  $\mathcal{E} = (\mathcal{X}, \mathcal{T}, \tau, E, \langle X_u \rangle_{0 \leq u \leq t})$  be a multi-objective optimization process with  $\preceq$  being a partial order on the multi-objective target space  $\mathcal{T}$ .

- A solution candidate  $x$  Pareto-dominates a solution candidate  $x'$  (assuming minimization) iff  $x \preceq x'$ .
- A solution candidate  $x$  is Pareto-optimal if there exists no other solution candidate  $x' \in \mathcal{X}$  so that  $x' \preceq x$ .
- The set of all Pareto-optimal solution candidates in  $\mathcal{X}$  is called the Pareto front of  $\mathcal{X}$  (w.r.t.  $\preceq$ ).

# Example: Multi-Objective Evolutionary Algorithm

# The Goal Class Hierarchy





## Goal Class 3: Goal Direction

*"I know which way it's getting better!"*

# Gradients of Target Functions

# Derivation

**Algorithm 5** (gradient descent). Let  $\mathcal{D} = (\mathcal{X}, \mathcal{T}, \tau, e, \langle x_u \rangle_{0 \leq u \leq t})$  be an optimization process. Let  $\mathcal{T}$  be continuous ( $\mathcal{T} = \mathbb{R}$ , e.g.) and let  $\tau' : \mathcal{X} \rightarrow \mathcal{T}$  be the first derivative of  $\tau$ . The process  $\mathcal{D}$  continues via gradient descent (with update rate  $\alpha \in \mathbb{R}^+$ ) if  $e$  is of the form

$$e(\langle x_u \rangle_{0 \leq u \leq t}, \tau) = x_{t+1} = x_t - \alpha \cdot \tau'(x_t).$$

The learning rate  $\alpha$  can also be given as a function, usually  $\alpha : \mathbb{N} \rightarrow \mathbb{R}$  so that  $e(\langle x_u \rangle_{0 \leq u \leq t}, \tau) = x_{t+1} = x_t - \alpha(t) \cdot \tau'(x_t)$ . If  $\tau$  is stochastic, this process is called stochastic gradient descent (SGD).

**Algorithm 5** (gradient descent). Let  $\mathcal{D} = (\mathcal{X}, \mathcal{T}, \tau, e, \langle x_u \rangle_{0 \leq u \leq t})$  be an optimization process. Let  $\mathcal{T}$  be continuous ( $\mathcal{T} = \mathbb{R}$ , e.g.) and let  $\tau' : \mathcal{X} \rightarrow \mathcal{T}$  be the first derivative of  $\tau$ . The process  $\mathcal{D}$  continues via gradient descent (with update rate  $\alpha \in \mathbb{R}^+$ ) if  $e$  is of the form

$$e(\langle x_u \rangle_{0 \leq u \leq t}, \tau) = x_{t+1} = x_t - \alpha \cdot \tau'(x_t).$$

The learning rate  $\alpha$  can also be given as a function, usually  $\alpha : \mathbb{N} \rightarrow \mathbb{R}$  so that  $e(\langle x_u \rangle_{0 \leq u \leq t}, \tau) = x_{t+1} = x_t - \alpha(t) \cdot \tau'(x_t)$ . If  $\tau$  is stochastic, this process is called stochastic gradient descent (SGD).

**Algorithm 6** (gradient descent (policy)). Let  $\pi_\theta$  be a policy  $\pi$  that depends on vector of continuous parameters  $\theta \in \Theta$  such that usually  $\Theta = \mathbb{R}^N$  for some  $N$ . Let  $\tau : \Theta \rightarrow \mathcal{T}$  be a target function on the parameters  $\theta$  of a policy  $\pi_\theta$ . Let  $\mathcal{T}$  be continuous ( $\mathcal{T} = \mathbb{R}$ , e.g.) and let  $\tau' : \Theta \rightarrow \mathcal{T}$  be the first derivative of  $\tau$ , i.e.,  $\tau'(\theta) = \frac{\partial \tau(\theta)}{\partial \theta}$ . If  $\mathcal{D} = (\Theta, \mathcal{T}, \tau, e, \langle x_u \rangle_{0 \leq u \leq t})$  is an optimization process that continues via gradient descent,  $\mathcal{D}$  is a process of policy optimization via gradient descent.

# Many Variants of Gradient Descent

source:

<https://optimization.cbe.cornell.edu/index.php?title=AdaGrad>

---

**Algorithm 1:** AdaGrad general algorithm

---

```
 $\eta$ : Stepsize ;  
 $f(x)$ : Stochastic objective function ;  
 $x_1$ : Initial parameter vector;  
for  $t = 1$  to  $T$  do  
    Evaluate  $f_t(x_t)$  ;  
    Get and save  $g_t$  ;  
     $G_t \leftarrow \sum_{\tau=1}^t g_\tau g_\tau^\top$  ;  
     $x_{t+1} \leftarrow x_t - \eta G_t^{-1/2} g_t$  ;  
end  
return  $x_t$ 
```

---

source: <https://arxiv.org/pdf/1412.6980.pdf%5D>

---

**Algorithm 1:** *Adam*, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation.  $g_t^2$  indicates the elementwise square  $g_t \odot g_t$ . Good default settings for the tested machine learning problems are  $\alpha = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 10^{-8}$ . All operations on vectors are element-wise. With  $\beta_1^t$  and  $\beta_2^t$  we denote  $\beta_1$  and  $\beta_2$  to the power  $t$ .

---

**Require:**  $\alpha$ : Stepsize

**Require:**  $\beta_1, \beta_2 \in [0, 1]$ : Exponential decay rates for the moment estimates

**Require:**  $f(\theta)$ : Stochastic objective function with parameters  $\theta$

**Require:**  $\theta_0$ : Initial parameter vector

$m_0 \leftarrow 0$  (Initialize 1<sup>st</sup> moment vector)

$v_0 \leftarrow 0$  (Initialize 2<sup>nd</sup> moment vector)

$t \leftarrow 0$  (Initialize timestep)

**while**  $\theta_t$  not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_\theta f_t(\theta_{t-1})$  (Get gradients w.r.t. stochastic objective at timestep  $t$ )

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$  (Update biased first moment estimate)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$  (Update biased second raw moment estimate)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$  (Compute bias-corrected first moment estimate)

$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$  (Compute bias-corrected second raw moment estimate)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$  (Update parameters)

**end while**

**return**  $\theta_t$  (Resulting parameters)

---