

- that determines the order in which S-polynomials are processed.
- match or beat the existing state-of-the-art.

nonzero ideal generated by polynomials f_1, \ldots, f_k .

SUMMARY **REINFORCEMENT LEARNING** Reinforcement learning problems can be phrased as the interaction of an agent and an environment. Buchberger's algorithm is the standard method for computing a Gröbner basis, and highly-tuned and optimized versions are a critical part of many computer algebra systems. 2. The efficiency of Buchberger's algorithm strongly depends on a choice of selection strategy Agent 3. By phrasing Buchberger's algorithm as a reinforcement learning problem and applying standard reinforcement learning techniques we can learn new selection strategies that can state reward S_t Gröbner Bases Environment Let $R = K[x_1, \ldots, x_n]$ be a polynomial ring over some field K and $I = \langle f_1, \ldots, f_k \rangle \subseteq R$ be a Given a monomial order, a Gröbner basis G of a I is a set of generators $\{g_1, g_2, \ldots, g_s\}$ of I The agent chooses actions and the environment processes actions and gives back the updated state and a such that any of the following equivalent conditions hold: reward. The agent wants to maximize its return, which is the amount of reward it gets in the long run. Chess CartPole f^{G} is unique for all $f \in R$ (iii) $\langle LT(g_1), LT(g_2), \dots, LT(g_s) \rangle = \langle LT(I) \rangle$ ◇ ≙ ≝ ☆ ≜ ◇ ! where LT(g) is the lead term of g with respect to the monomial order and $f^G \rightarrow r$ is the remainder under polynomial long division of f by the polynomials in G. BUCHBERGER'S ALGORITHM **Theorem (Buchberger's Criterion)**: Let $G = \{g_1, g_2, \dots, g_s\}$ generate some ideal *I*. If $S(g_i, g_i)^G \rightarrow 0$ for all pairs g_i, g_i , where S_t = positions of all pieces $S_t = \text{cart/pole position/velocity}$ A_t = push the cart left or right $A_t =$ a valid move $R_t = 1$ while pole is upright $R_t = 1$ if you win at t, otherwise 0 POLICIES AND TRAJECTORIES **input** a set of polynomials $\{f_1, \ldots, f_k\}$ **output** a Gröbner basis G of $I = \langle f_1, \ldots, f_k \rangle$ **procedure** BUCHBERGER($\{f_1, \ldots, f_k\}$) A policy π is a function $\pi : \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ given by ▷ the current basis $G \leftarrow \{f_1, \ldots, f_k\}$ $\pi(a|s) = \Pr(A_t = a|S_t = s)$ $P \leftarrow \{(f_i, f_j) \mid 1 \leq i < j \leq k\}$ ▷ the remaining pairs which maps state-action pairs to the probability of choosing the given action in the given state. while |P| > 0 do Policies are often viewed as functions that take in a state and return a probability distribution on actions. An $(f_i, f_i) \leftarrow \operatorname{select}(P)$ agent follows a policy by applying the policy to its current state and sampling from the returned probability $P \leftarrow P \setminus \{(f_i, f_j)\}$ distribution to choose the next action. $r \leftarrow S(f_i, f_i)^G$ A *trajectory* or *rollout* τ of a policy π is a series of states, actions, and rewards **if** *r* ≠ 0 **then** $G \leftarrow G \cup \{r\}$ $\tau = (S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \dots, R_T, S_T)$ $P \leftarrow P \cup \{(f, r) : f \in G\}$ obtained by following the policy π one time through the environment, and the *return* of a trajectory is the end if sum of rewards along the trajectory. end while return G Given an environment, the goal of reinforcement learning is to find a policy π that maximizes end procedure $\mathop{\mathbb{E}}_{\tau\sim\pi}$ SELECTION STRATEGIES IN BUCHBERGER'S ALGORITHM which is the expected return along trajectories τ obtained by following π . The implementation of select does not affect correctness of Buchberger's algorithm, but it is critical for efficiency. In general, good selection strategies pick "small" pairs first. POLICY GRADIENT First: among the pairs with minimal *j*, pick the pair with smallest *i* ▶ Degree: pick the pair with smallest degree of $Icm(LT(f_i), LT(f_i))$ Suppose π_{θ} is a parametrized policy that is differentiable with respect to its parameters θ . Then the ▶ Normal: pick the pair with smallest $lcm(LT(f_i), LT(f_i))$ in the monomial order expected return Sugar: pick the pair with smallest sugar degree of $Icm(LT(f_i), LT(f_i))$ $J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left| \sum_{t=1}^{T} R_{t} \right|$ Pair Reductions in Buchberger's Algorithm per Strategy ce has gradient $\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(A_t | S_t) \sum_{t'=t+1}^{T} R_{t'} \right].$ This expectation is a quantity we can sample by interacting with the environment. By starting with any set of parameters θ_1 and updating by gradient ascent steps $\theta_{k} = \theta_{k-1} + \alpha \cdot \nabla_{\theta} J(\theta_{k})$ for some small learning rate α , we can incrementally improve the policy. Intuitively, we should increase the

(i)
$$f^G \to 0 \iff f \in I$$

the log probability of choosing that action again.

$$S(g_i, g_j) = \frac{\operatorname{lcm}(\operatorname{LT}(g_i), \operatorname{LT}(g_j))}{\operatorname{LT}(g_i)} g_i - \frac{\operatorname{lcm}(\operatorname{LT}(g_i), \operatorname{LT}(g_j))}{\operatorname{LT}(g_i)}$$

is the S-polynomial of g_i and g_i , then G is a Gröbner basis of I.

Algorithm 1 Buchberger's Algorithm

r an rioddollorio in Edolloorgor o rigoritini por olidlogy								
example	First	Normal	Sugar	Random	Last	Strange	Spic	
cyclic4	11	11	11	14	21	23	23	
reimer3	25	23	25	25	25	28	28	
katsura5	28	28	28	44	76	86	86	
eco6	67	61	64	97	149	295	295	
noon4	71	71	71	100	103	375	375	
cyclic6	366	620	343	793	-	-	-	
katsura7	164	164	164	285	-	-	-	
katsura4-lex	25	46	19	29	44	30	59	
eco5-lex	30	22	26	28	91	32	97	
cyclic5-lex	104	1602	108	-	-	-	-	
	1			1	1			

Reinforcement Learning in Buchberger's Algorithm Dylan Peifer **Cornell University**





