

Learning Selection Strategies in Buchberger's Algorithm

Dylan Peifer

Cornell University

23 October 2020

Outline

The efficiency of Buchberger's algorithm, a fundamental tool in computer algebra, strongly depends on a choice of selection strategy. By phrasing Buchberger's algorithm as a reinforcement learning problem and applying standard reinforcement learning techniques we can learn new selection strategies that beat the existing state-of-the-art.

1. Gröbner Bases and Buchberger's Algorithm
2. Reinforcement Learning and Policy Gradient
3. Results

1. Gröbner Bases and Buchberger's Algorithm

$R = K[x_1, \dots, x_n]$ a polynomial ring over some field K

$I = \langle f_1, \dots, f_k \rangle \subseteq R$ an ideal generated by $f_1, \dots, f_k \in R$

$R = K[x_1, \dots, x_n]$ a polynomial ring over some field K

$I = \langle f_1, \dots, f_k \rangle \subseteq R$ an ideal generated by $f_1, \dots, f_k \in R$

Example

$$\begin{aligned} R &= \mathbb{Q}[x, y] \\ &= \{ \text{polynomials in } x \text{ and } y \text{ with rational coefficients} \} \end{aligned}$$

$$\begin{aligned} I &= \langle x^2 - y^3, xy^2 + x \rangle \\ &= \{ a(x^2 - y^3) + b(xy^2 + x) : a, b \in R \} \end{aligned}$$

$R = K[x_1, \dots, x_n]$ a polynomial ring over some field K

$I = \langle f_1, \dots, f_k \rangle \subseteq R$ an ideal generated by $f_1, \dots, f_k \in R$

Example

$$\begin{aligned} R &= \mathbb{Q}[x, y] \\ &= \{\text{polynomials in } x \text{ and } y \text{ with rational coefficients}\} \end{aligned}$$

$$\begin{aligned} I &= \langle x^2 - y^3, xy^2 + x \rangle \\ &= \{a(x^2 - y^3) + b(xy^2 + x) : a, b \in R\} \end{aligned}$$

Question

In the above example, is $x^5 + x$ an element of I ?

Question

Consider the ideal $I = \langle x^2 + x - 2 \rangle$ in the ring $\mathbb{Q}[x]$. Is $x^3 + 3x^2 + 5x + 4$ an element of I ?

Question

Consider the ideal $I = \langle x^2 + x - 2 \rangle$ in the ring $\mathbb{Q}[x]$. Is $x^3 + 3x^2 + 5x + 4$ an element of I ?

$$\begin{array}{r}
 x^2 + x - 2 \quad) \quad \begin{array}{r}
 x^3 + 3x^2 + 5x + 4 \\
 - (x^3 + x^2 - 2x) \\
 \hline
 2x^2 + 7x + 4 \\
 - (2x^2 + 2x - 4) \\
 \hline
 5x + 8
 \end{array}
 \end{array}$$

Question

Consider the ideal $I = \langle x^2 + x - 2 \rangle$ in the ring $\mathbb{Q}[x]$. Is $x^3 + 3x^2 + 5x + 4$ an element of I ?

$$\begin{array}{r} x^2 + x - 2 \overline{) \begin{array}{r} x^3 + 3x^2 + 5x + 4 \\ - (x^3 + x^2 - 2x) \\ \hline 2x^2 + 7x + 4 \\ - (2x^2 + 2x - 4) \\ \hline 5x + 8 \end{array}} \end{array}$$

$$x^3 + 3x^2 + 5x + 4 = (x + 2)(x^2 + x - 2) + (5x + 8)$$

Question

Consider the ideal $I = \langle x^2 + x - 2 \rangle$ in the ring $\mathbb{Q}[x]$. Is $x^3 + 3x^2 + 5x + 4$ an element of I ?

$$\begin{array}{r} x^2 + x - 2 \overline{) \quad \begin{array}{r} x^3 + 3x^2 + 5x + 4 \\ - (x^3 + x^2 - 2x) \\ \hline 2x^2 + 7x + 4 \\ - (2x^2 + 2x - 4) \\ \hline 5x + 8 \end{array}} \end{array}$$

$$x^3 + 3x^2 + 5x + 4 = (x + 2)(x^2 + x - 2) + (5x + 8)$$

$$\Rightarrow \boxed{x^3 + 3x^2 + 5x + 4 \notin \langle x^2 + x - 2 \rangle}$$

Definition

Let x^α denote an arbitrary monomial where α is the vector of exponents. A **monomial order** on $R = k[x_1, \dots, x_n]$ is a relation $>$ on the monomials of R such that

1. $>$ is a total ordering
2. $>$ is a well-ordering
3. if $x^\alpha > x^\beta$ then $x^\gamma x^\alpha > x^\gamma x^\beta$ for any x^γ (i.e., $>$ respects multiplication).

Definition

Let x^α denote an arbitrary monomial where α is the vector of exponents. A **monomial order** on $R = k[x_1, \dots, x_n]$ is a relation $>$ on the monomials of R such that

1. $>$ is a total ordering
2. $>$ is a well-ordering
3. if $x^\alpha > x^\beta$ then $x^\gamma x^\alpha > x^\gamma x^\beta$ for any x^γ (i.e., $>$ respects multiplication).

Example

Lexicographic order (lex) is defined by $x^\alpha > x^\beta$ if the leftmost nonzero component of $\alpha - \beta$ is positive. For example, $x > y > z$, $xy > y^4$, and $xz > y^2$.

Divide $x^5 + x$ by the generators $x^2 - y^3$ and $xy^2 + x$

$$\begin{array}{r}
 \begin{array}{cc} x^2 & - \\ xy^2 & + \end{array} \begin{array}{cc} y^3 & \\ x & \end{array} & \begin{array}{l} q_1 : x^3 - xy \\ q_2 : x^2y - y^2 + 1 \end{array} \\
 \hline
 - \begin{array}{cc} x^5 & + \\ (x^5 & - \end{array} \begin{array}{cc} x & \\ x^3y^3 & \end{array} \\
 \hline
 \begin{array}{cc} x^3y^3 & + \\ (x^3y^3 & + \end{array} \begin{array}{cc} x & \\ x^3y & \end{array} \\
 \hline
 - \begin{array}{cc} x^3y & + \\ (-x^3y & + \end{array} \begin{array}{cc} x & \\ xy^4 & \end{array} \\
 \hline
 - \begin{array}{cc} xy^4 & + \\ (-xy^4 & - \end{array} \begin{array}{cc} x & \\ xy^2 & \end{array} \\
 \hline
 \begin{array}{cc} xy^2 & + \\ (xy^2 & + \end{array} \begin{array}{cc} x & \\ x & \end{array} \\
 \hline
 0
 \end{array}$$

Divide $x^5 + x$ by the generators $x^2 - y^3$ and $xy^2 + x$

$$\begin{array}{r}
 \begin{array}{c} x^2 \\ xy^2 \end{array} \quad \begin{array}{c} - \\ + \end{array} \quad \begin{array}{c} y^3 \\ x \end{array} \\
 \hline
 \begin{array}{r}
 q_1 : \quad x^3 \quad - \quad xy \\
 q_2 : \quad x^2y \quad - \quad y^2 \quad + \quad 1 \\
 \hline
 \quad \quad x^5 \quad + \quad x \\
 \quad \quad (x^5 \quad - \quad x^3y^3) \\
 \hline
 \quad \quad \quad x^3y^3 \quad + \quad x \\
 \quad \quad \quad (x^3y^3 \quad + \quad x^3y) \\
 \hline
 \quad \quad \quad \quad -x^3y \quad + \quad x \\
 \quad \quad \quad \quad (-x^3y \quad + \quad xy^4) \\
 \hline
 \quad \quad \quad \quad \quad -xy^4 \quad + \quad x \\
 \quad \quad \quad \quad \quad (-xy^4 \quad - \quad xy^2) \\
 \hline
 \quad \quad \quad \quad \quad \quad xy^2 \quad + \quad x \\
 \quad \quad \quad \quad \quad \quad (-xy^2 \quad + \quad x) \\
 \hline
 \quad \quad \quad \quad \quad \quad \quad 0
 \end{array}
 \end{array}$$

$$x^5 + x = (x^3 - xy)(x^2 - y^3) + (x^2y - y^2 + 1)(xy^2 + x) + 0$$

Divide $x^5 + x$ by the generators $x^2 - y^3$ and $xy^2 + x$

$$\begin{array}{r}
 \begin{array}{c} x^2 \quad - \quad y^3 \\ xy^2 \quad + \quad x \end{array} \quad \begin{array}{r} q_1 : \quad x^3 \quad - \quad xy \\ q_2 : \quad x^2y \quad - \quad y^2 \quad + \quad 1 \\ \hline - \quad x^5 \quad + \quad x^3y^3 \\ \hline \quad \quad - \quad x^3y^3 \quad + \quad x^3y^3 \\ \hline \quad \quad \quad - \quad x^3y^3 \quad + \quad x^3y^3 \\ \hline \quad \quad \quad \quad - \quad x^3y^3 \quad + \quad xy^4 \\ \hline \quad \quad \quad \quad \quad - \quad xy^4 \quad + \quad xy^2 \\ \hline \quad \quad \quad \quad \quad \quad - \quad xy^2 \quad + \quad x \\ \hline \quad \quad \quad \quad \quad \quad \quad - \quad xy^2 \quad + \quad x \\ \hline \quad \quad \quad \quad \quad \quad \quad \quad 0 \end{array}
 \end{array}$$

$$x^5 + x = (x^3 - xy)(x^2 - y^3) + (x^2y - y^2 + 1)(xy^2 + x) + 0$$

$$\Rightarrow \boxed{x^5 + x \in \langle x^2 - y^3, xy^2 + x \rangle}$$

Definition

When F is set of polynomials and dividing h by the $f_i \in F$ using the division algorithm leads to the remainder r we write $h^F \rightarrow r$ or say h reduces to r .

Definition

When F is set of polynomials and dividing h by the $f_i \in F$ using the division algorithm leads to the remainder r we write $h^F \rightarrow r$ or say h reduces to r .

Lemma

If $h^F \rightarrow 0$ then h is in the ideal generated by F .

Definition

When F is set of polynomials and dividing h by the $f_i \in F$ using the division algorithm leads to the remainder r we write $h^F \rightarrow r$ or say h reduces to r .

Lemma

If $h^F \rightarrow 0$ then h is in the ideal generated by F .

Unfortunately, the converse is false.

Example

Using the same ideal $I = \langle x^2 - y^3, xy^2 + x \rangle$, note that

$$y^2(x^2 - y^3) - x(xy^2 + x) = -x^2 - y^5 \in I$$

However, multivariate division produces the nonzero remainder $-y^5 - y^3$.

Definition

Given a monomial order, a **Gröbner basis** G of a nonzero ideal I is a set of generators $\{g_1, g_2, \dots, g_s\}$ of I such that any of the following equivalent conditions hold:

- (i) $f^G \rightarrow 0 \iff f \in I$
- (ii) f^G is unique for all $f \in R$
- (iii) $\langle \text{LT}(g_1), \text{LT}(g_2), \dots, \text{LT}(g_s) \rangle = \langle \text{LT}(I) \rangle$

where $\text{LT}(f)$ is the **leading term** of f and $\langle \text{LT}(I) \rangle = \langle \text{LT}(f) \mid f \in I \rangle$ is the ideal generated by all leading terms of I .

Definition

Given a monomial order, a **Gröbner basis** G of a nonzero ideal I is a set of generators $\{g_1, g_2, \dots, g_s\}$ of I such that any of the following equivalent conditions hold:

- (i) $f^G \rightarrow 0 \iff f \in I$
- (ii) f^G is unique for all $f \in R$
- (iii) $\langle \text{LT}(g_1), \text{LT}(g_2), \dots, \text{LT}(g_s) \rangle = \langle \text{LT}(I) \rangle$

where $\text{LT}(f)$ is the **leading term** of f and $\langle \text{LT}(I) \rangle = \langle \text{LT}(f) \mid f \in I \rangle$ is the ideal generated by all leading terms of I .

Example

Using the same ideal $I = \langle x^2 - y^3, xy^2 + x \rangle$, the set $\{x^2 - y^3, xy^2 + x\}$ is **not** a Gröbner basis of I .

Definition

Let $S(f, g) = \frac{x^\gamma}{\text{LT}(f)}f - \frac{x^\gamma}{\text{LT}(g)}g$ where x^γ is the least common multiple of the leading monomials of f and g . This is the *s-polynomial* of f and g , where s stands for subtraction or syzygy.

Definition

Let $S(f, g) = \frac{x^\gamma}{\text{LT}(f)}f - \frac{x^\gamma}{\text{LT}(g)}g$ where x^γ is the least common multiple of the leading monomials of f and g . This is the *s-polynomial* of f and g , where s stands for subtraction or syzygy.

Example

$$\begin{aligned} S(x^2 - y^3, xy^2 + x) &= \frac{x^2y^2}{x^2}(x^2 - y^3) - \frac{x^2y^2}{xy^2}(xy^2 + x) \\ &= y^2(x^2 - y^3) - x(xy^2 + x) \\ &= -x^2 - y^5 \end{aligned}$$

Definition

Let $S(f, g) = \frac{x^\gamma}{\text{LT}(f)}f - \frac{x^\gamma}{\text{LT}(g)}g$ where x^γ is the least common multiple of the leading monomials of f and g . This is the *s-polynomial* of f and g , where s stands for subtraction or syzygy.

Example

$$\begin{aligned} S(x^2 - y^3, xy^2 + x) &= \frac{x^2y^2}{x^2}(x^2 - y^3) - \frac{x^2y^2}{xy^2}(xy^2 + x) \\ &= y^2(x^2 - y^3) - x(xy^2 + x) \\ &= -x^2 - y^5 \end{aligned}$$

Theorem (Buchberger's Criterion)

Let $G = \{g_1, g_2, \dots, g_s\}$ generate the ideal I . If $S(g_i, g_j)^G \rightarrow 0$ for all pairs g_i, g_j then G is a Gröbner basis of I .

Algorithm Buchberger's Algorithm

input a set of polynomials $\{f_1, \dots, f_k\}$

output a Gröbner basis G of $I = \langle f_1, \dots, f_k \rangle$

procedure BUCHBERGER($\{f_1, \dots, f_k\}$)

$$G \leftarrow \{f_1, \dots, f_k\}$$

▷ the current basis

$$P \leftarrow \{(f_i, f_j) \mid 1 \leq i < j \leq k\}$$

- ▷ the remaining pairs

while $|P| > 0$ **do**
$$(f_i, f_j) \leftarrow \text{select}(P)$$
$$P \leftarrow P \setminus \{(f_i, f_j)\}$$
$$r \leftarrow S(f_i, f_j)^G$$

if $r \neq 0$ then

$$P \leftarrow P \cup \{(f, r) : f \in G\}$$
$$G \leftarrow G \cup \{r\}$$

end if

end while

```

return G

```

end procedure

Example

$$I = \langle x^2 - y^3, xy^2 + x \rangle$$

Example

$$I = \langle x^2 - y^3, xy^2 + x \rangle$$

initialize G to $\{x^2 - y^3, xy^2 + x\}$

initialize P to $\{(x^2 - y^3, xy^2 + x)\}$

Example

$$I = \langle x^2 - y^3, xy^2 + x \rangle$$

initialize G to $\{x^2 - y^3, xy^2 + x\}$

initialize P to $\{(x^2 - y^3, xy^2 + x)\}$

select $(x^2 - y^3, xy^2 + x)$ and compute $S(x^2 - y^3, xy^2 + x)^G \rightarrow -y^5 - y^3$

update G to $\{x^2 - y^3, xy^2 + x, -y^5 - y^3\}$

update P to $\{(x^2 - y^3, -y^5 - y^3), (xy^2 + x, -y^5 - y^3)\}$

Example

$$I = \langle x^2 - y^3, xy^2 + x \rangle$$

initialize G to $\{x^2 - y^3, xy^2 + x\}$

initialize P to $\{(x^2 - y^3, xy^2 + x)\}$

select $(x^2 - y^3, xy^2 + x)$ and compute $S(x^2 - y^3, xy^2 + x)^G \rightarrow -y^5 - y^3$

update G to $\{x^2 - y^3, xy^2 + x, -y^5 - y^3\}$

update P to $\{(x^2 - y^3, -y^5 - y^3), (xy^2 + x, -y^5 - y^3)\}$

select $(x^2 - y^3, -y^5 - y^3)$ and compute $S(x^2 - y^3, -y^5 - y^3)^G \rightarrow 0$

Example

$$I = \langle x^2 - y^3, xy^2 + x \rangle$$

initialize G to $\{x^2 - y^3, xy^2 + x\}$

initialize P to $\{(x^2 - y^3, xy^2 + x)\}$

select $(x^2 - y^3, xy^2 + x)$ and compute $S(x^2 - y^3, xy^2 + x)^G \rightarrow -y^5 - y^3$

update G to $\{x^2 - y^3, xy^2 + x, -y^5 - y^3\}$

update P to $\{(x^2 - y^3, -y^5 - y^3), (xy^2 + x, -y^5 - y^3)\}$

select $(x^2 - y^3, -y^5 - y^3)$ and compute $S(x^2 - y^3, -y^5 - y^3)^G \rightarrow 0$

select $(xy^2 + x, -y^5 - y^3)$ and compute $S(xy^2 + x, -y^5 - y^3)^G \rightarrow 0$

Example

$$I = \langle x^2 - y^3, xy^2 + x \rangle$$

initialize G to $\{x^2 - y^3, xy^2 + x\}$

initialize P to $\{(x^2 - y^3, xy^2 + x)\}$

select $(x^2 - y^3, xy^2 + x)$ and compute $S(x^2 - y^3, xy^2 + x)^G \rightarrow -y^5 - y^3$

update G to $\{x^2 - y^3, xy^2 + x, -y^5 - y^3\}$

update P to $\{(x^2 - y^3, -y^5 - y^3), (xy^2 + x, -y^5 - y^3)\}$

select $(x^2 - y^3, -y^5 - y^3)$ and compute $S(x^2 - y^3, -y^5 - y^3)^G \rightarrow 0$

select $(xy^2 + x, -y^5 - y^3)$ and compute $S(xy^2 + x, -y^5 - y^3)^G \rightarrow 0$

return $G = \{x^2 - y^3, xy^2 + x, -y^5 - y^3\}$

Algorithm Buchberger's Algorithm

input a set of polynomials $\{f_1, \dots, f_k\}$

output a Gröbner basis G of $I = \langle f_1, \dots, f_k \rangle$

procedure BUCHBERGER($\{f_1, \dots, f_k\}$)

$$G \leftarrow \{f_1, \dots, f_k\}$$

▷ the current basis

$$P \leftarrow \{(f_i, f_j) \mid 1 \leq i < j \leq k\}$$

- ▷ the remaining pairs

while $|P| > 0$ **do**
$$(f_i, f_j) \leftarrow \text{select}(P)$$
$$P \leftarrow P \setminus \{(f_i, f_j)\}$$
$$r \leftarrow S(f_i, f_j)^G$$

if $r \neq 0$ then

$$P \leftarrow P \cup \{(f, r) : f \in G\}$$
$$G \leftarrow G \cup \{r\}$$

end if

end while

```

return G

```

end procedure

In general, we should select “small” pairs (f_i, f_j) first.

In general, we should select “small” pairs (f_i, f_j) first.

- ▶ First:
among the pairs with minimal j , pick the pair with smallest i
- ▶ Degree:
pick the pair with smallest degree of $\text{lcm}(\text{LT}(f_i), \text{LT}(f_j))$
- ▶ Normal:
pick the pair with smallest $\text{lcm}(\text{LT}(f_i), \text{LT}(f_j))$ in the monomial order
- ▶ Sugar:
pick the pair with smallest sugar degree of $\text{lcm}(\text{LT}(f_i), \text{LT}(f_j))$, which is the degree it would have had if we had homogenized at the beginning

The number of pair reductions performed is a rough estimate of how much time was spent. Smaller numbers are better.

example	First	Degree	Normal	Sugar	Random
cyclic6	371	655	620	343	793
cyclic7	2217	5664	5781	2070	-
katsura7	164	164	164	164	285
eco6	67	72	61	64	97
reimer5	552	212	211	301	-
noon4	71	71	71	71	100
cyclic5 (lex)	112	132	1602	108	-
katsura5 (lex)	231	1631	769	67	-
eco5 (lex)	30	34	22	26	28
eco6 (lex)	104	147	96	68	175

Summary

- ▶ A Gröbner basis of an ideal in a polynomial ring is a special generating set that is useful for many computational problems.
- ▶ Buchberger's algorithm produces a Gröbner basis from any initial generating set of an ideal by repeatedly choosing pairs (f_i, f_j) of the current generating set and adding the reduction of the s-polynomial of f_i and f_j to the generating set if it is not zero.
- ▶ The selection strategy used to pick which pair to choose next can make a big difference in the efficiency of Buchberger's algorithm.

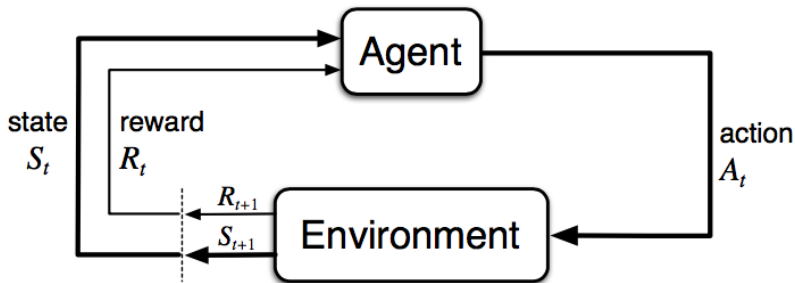
2. Reinforcement Learning and Policy Gradient

Reinforcement learning tries to understand and optimize goal-directed behavior driven by interaction with the world.

Reinforcement learning tries to understand and optimize goal-directed behavior driven by interaction with the world.

- ▶ playing games (backgammon, chess, Go, StarCraft, ...)
- ▶ flying a helicopter or driving a car
- ▶ controlling a power station or data center
- ▶ managing a portfolio of stocks or other financial assets
- ▶ allocating resources to research projects

Reinforcement learning problems can be phrased as the interaction of an agent and an environment.



The agent chooses actions and the environment processes actions and gives back the updated state and a reward. The agent wants to maximize its return, which is the amount of reward it gets in the long run.

Definition

A *Markov Decision Process (MDP)* is a collection of states \mathcal{S} and actions \mathcal{A} with transition dynamics given by

$$p : \mathcal{S} \times \mathbb{R} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$$

where

$$p(s', r | s, a) = \Pr[S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a]$$

returns the probability that the next state is s' and the next reward is r given that the current state is s and the chosen action is a .

Definition

A *Markov Decision Process (MDP)* is a collection of states \mathcal{S} and actions \mathcal{A} with transition dynamics given by

$$p : \mathcal{S} \times \mathbb{R} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$$

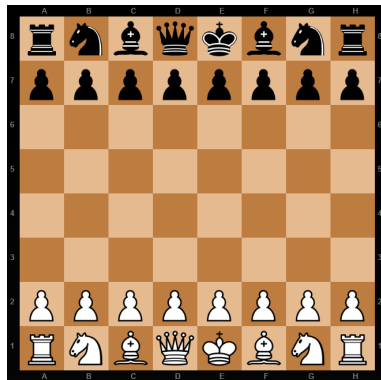
where

$$p(s', r | s, a) = \Pr[S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a]$$

returns the probability that the next state is s' and the next reward is r given that the current state is s and the chosen action is a .

An environment implements an MDP by computing $p(\cdot, \cdot | s, a)$ for the current state s and action a provided by the agent and then sampling from the resulting distribution to return a new state s' and reward r .

Chess



State: the positions of all pieces on the board

Action: a valid move of one of your pieces

Reward: 1 if you win immediately after the transition, otherwise 0

Definition

A *policy* π is a function

$$\pi : \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$$

where

$$\pi(a|s) = \Pr(A_t = a | S_t = s)$$

returns the probability that the next action is a given that the current state is s .

Definition

A *policy* π is a function

$$\pi : \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$$

where

$$\pi(a|s) = Pr(A_t = a | S_t = s)$$

returns the probability that the next action is a given that the current state is s .

An agent follows a policy by computing $\pi(\cdot|s)$ for the current state s and sampling from the resulting probability distribution to choose the next action.

Definition

A *trajectory*, *episode*, or *rollout* τ of a policy π is a series of states, actions, and rewards $(S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \dots, R_T, S_T)$ obtained by following the policy π one time through the environment.

Definition

The *return* of a trajectory is the sum of rewards

$$\sum_{t=1}^T R_t$$

along the trajectory.

The Reinforcement Learning Problem

Given an MDP, determine a policy π that maximizes the expected return

$$\mathbb{E}_{\tau \sim \pi} \left[\sum_{t=1}^T R_t \right]$$

over full trajectories sampled by following the policy π .

The Reinforcement Learning Problem

Given an MDP, determine a policy π that maximizes the expected return

$$\mathbb{E}_{\tau \sim \pi} \left[\sum_{t=1}^T R_t \right]$$

over full trajectories sampled by following the policy π .

If we know the exact transition dynamics of the MDP this is a **planning** problem. In the full **learning** problem the dynamics are either unknown or infeasible to compute. All we can do is sample from the environment.

Consider a parametrized policy function π_θ which maps states to probability distributions on actions. The expected return is now a function

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=1}^T R_t \right]$$

of the parameters θ of the policy.

Consider a parametrized policy function π_θ which maps states to probability distributions on actions. The expected return is now a function

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=1}^T R_t \right]$$

of the parameters θ of the policy.

Starting from any value of the parameters θ_1 , we can improve the policy by repeatedly moving the parameters in the direction of $\nabla_\theta J(\theta)$

$$\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\theta)|_{\theta_k}$$

where α is some small learning rate.

Theorem (Policy Gradient Theorem)

Suppose π_θ is a parametrized policy that is differentiable with respect to its parameters θ . Then the gradient of

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=1}^T R_t \right]$$

is

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(A_t | S_t) \sum_{t'=t+1}^T R_{t'} \right].$$

Theorem (Policy Gradient Theorem)

Suppose π_θ is a parametrized policy that is differentiable with respect to its parameters θ . Then the gradient of

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=1}^T R_t \right]$$

is

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(A_t | S_t) \sum_{t'=t+1}^T R_{t'} \right].$$

Intuitively, we should increase the probability of taking the action we chose proportional to the future reward we received and the derivative of the log probability of choosing that action again.

Summary

- ▶ Reinforcement learning can be phrased as the interaction of an agent and an environment, where an agent picks actions and is trying to maximize the total reward it receives from the environment over a full trajectory.
- ▶ A policy is a function that takes in a state and returns a probability distribution on actions.
- ▶ Policy gradient methods improve a parametrized policy by moving the parameters in the direction of the gradient of expected return.

3. Results

Algorithm Buchberger's Algorithm

input a set of polynomials $\{f_1, \dots, f_k\}$

output a Gröbner basis G of $I = \langle f_1, \dots, f_k \rangle$

procedure BUCHBERGER($\{f_1, \dots, f_k\}$)

$G \leftarrow \{f_1, \dots, f_k\}$

▷ the current basis

$P \leftarrow \{(f_i, f_j) \mid 1 \leq i < j \leq k\}$

▷ the remaining pairs

while $|P| > 0$ **do**

$(f_i, f_j) \leftarrow \text{select}(P)$

$P \leftarrow P \setminus \{(f_i, f_j)\}$

$r \leftarrow S(f_i, f_j)^G$

if $r \neq 0$ **then**

$P \leftarrow P \cup \{(f, r) : f \in G\}$

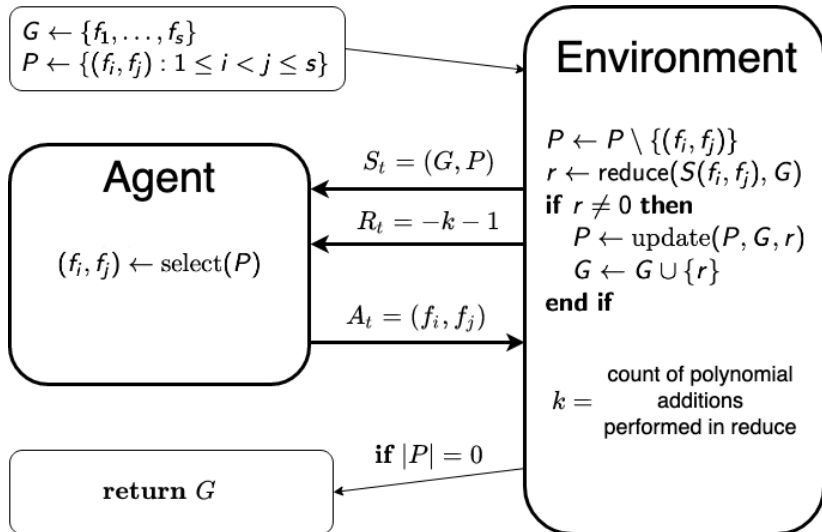
$G \leftarrow G \cup \{r\}$

end if

end while

return G

end procedure



Choosing a Distribution of Ideals

Starting generators are **binomials** with no constant terms in 3 variables and a fixed maximum degree.

Example

$$\{x^3z + y^2, \quad x^2z^2 - xyz, \quad x^2y - 3z\}$$

Choosing a Distribution of Ideals

Starting generators are **binomials** with no constant terms in 3 variables and a fixed maximum degree.

Example

$$\{x^3z + y^2, \quad x^2z^2 - xyz, \quad x^2y - 3z\}$$

- ▶ We avoid uninteresting generic behavior.
- ▶ All new generators are also binomial.
- ▶ Some of the hardest known examples are binomial ideals.
- ▶ By adjusting the degree and number of initial generators, we can adjust the difficulty of the problem.

Expressing the State to the Model

The state (G, P) is mapped to a $|P| \times 12$ matrix with each row given by the

$$(2 \text{ binomials})(2 \text{ terms})(3 \text{ variables}) = 12 \text{ exponents}$$

involved in each pair.

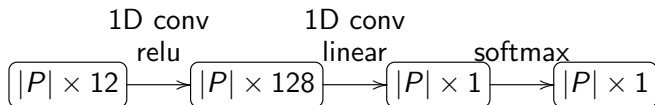
Expressing the State to the Model

The state (G, P) is mapped to a $|P| \times 12$ matrix with each row given by the

$$(2 \text{ binomials})(2 \text{ terms})(3 \text{ variables}) = 12 \text{ exponents}$$

involved in each pair.

This matrix is passed into a policy network



and a value model which computes the future return from following Degree selection.

The network weights are initialized randomly. Training then proceeds through epochs. In each epoch:

The network weights are initialized randomly. Training then proceeds through epochs. In each epoch:

1. Perform 100 rollouts using the current policy network.

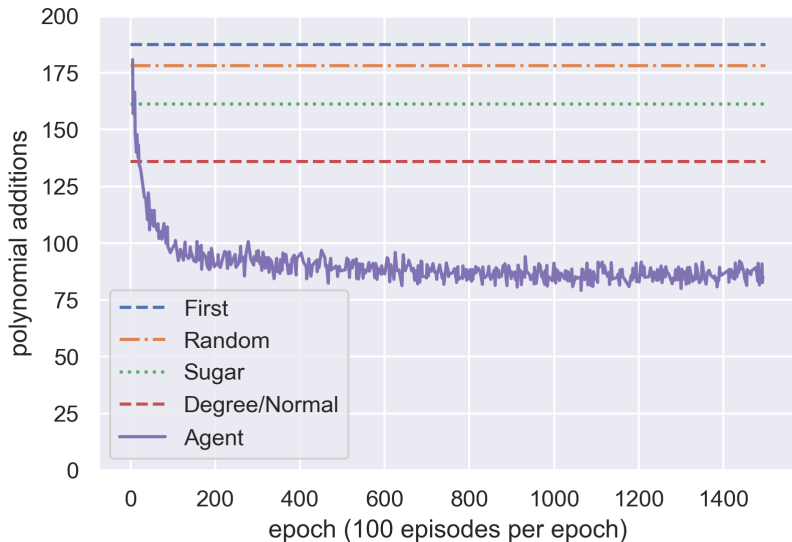
The network weights are initialized randomly. Training then proceeds through epochs. In each epoch:

1. Perform 100 rollouts using the current policy network.
2. Compute future rewards for each action on each trajectory, use generalized advantage estimation with a Degree model baseline, and normalize these scores across the epoch.

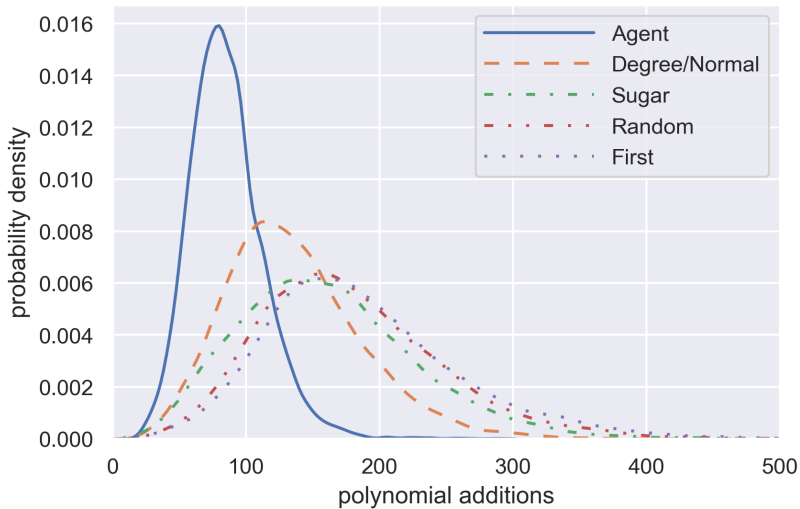
The network weights are initialized randomly. Training then proceeds through epochs. In each epoch:

1. Perform 100 rollouts using the current policy network.
2. Compute future rewards for each action on each trajectory, use generalized advantage estimation with a Degree model baseline, and normalize these scores across the epoch.
3. Update the policy network using gradient ascent and the policy gradient theorem.

Training



Testing



s	DIST	RANDOM	NORMAL	AGENT	IMPROVEMENT
10	W	178.[68.3]	136.[51.2]	85.6[27.3]	37% [46%]
4	W	203.[97.8]	160.[66.6]	101.[44.9]	37% [30%]
10	U	318.[103.]	198.[57.1]	141.[42.8]	28% [23%]
4	U	303.[122.]	194.[70.0]	151.[56.4]	22% [19%]

Agent performance in 3 variables and degree 20. Each line is a unique agent trained on the given distribution. Performance is mean[stddev] of polynomial additions on 10000 random samples.

Summary

- ▶ Pair selection, a key choice in Buchberger's algorithm, can be expressed as a reinforcement learning problem.
- ▶ Passing the state to a neural network is challenging since the state is unbounded in several directions, and training and testing requires choosing some distribution of ideals.
- ▶ In several distributions of random binomial ideals, our trained model outperforms state-of-the-art human-designed selection strategies by 20% to 40%.

Dylan Peifer
Michael Stillman
Daniel Halpern-Leistner

djp282@cornell.edu
mes15@cornell.edu
daniel.hl@cornell.edu

[1] Learning selection strategies in Buchberger's algorithm. In *Proceedings of the 37th International Conference on Machine Learning (ICML 2020)*.

[2] <https://github.com/dylanpeifer/deepgroebner>