Optimizing non-differentiable metrics

And applications in autonomous driving

Plan for Today

Jens Beißwenger

Andreas Geiger

2



Bernhard Jaeger, Andreas Geiger: "An Invitation to Deep Reinforcement Learning", Foundations and Trends® in Optimization 2024











Article Faster sorting algorithms discovered using deep reinforcement learning

Article							
Magnetic conf	troloftokamakr	plasmas					
thagheete com	· · ·	·					
through deep	reinforcement le	earning					
		5					
https://doi.org/10.1038/s41586-021-0430	Article	e					
Received: 14 July 2021	stor matrix multiplication						
Accepted: 1 December 2021	Discovering raster matrix multiplication						
Published online: 16 February 2022	itshed online: 16 February 2022 algorithms with reinforcement learning						
Open access	8	8					
Check for updates							
	https://doi.org/10.1038/s41586.022.05172.4	Albussein Fawzi ^{1/2}					
	Received: 2 October 2021	Bernardino Romera-Paredes ¹² , Mohammadamin Barekatain ¹ , Alexander Novikov ¹ ,					
	Accepted: 2 August 2022	Francisco J. R. Ruiz', Julian Schrittwieser', Grzegorz Swirszcz', David Silver', Demis Hassabis & Pushmeet Kohli ¹					
	Published online: 5 October 2022						
	Open access	Improving the efficiency of algorithms for fundamental computations can have a					
	Check for updates	widespread impact, as it can affect the overall speed of a large amount of computatic Matrix multiplication is one such primitive task, occurring in many systems – from neural networks to scientific computing routines. The automatic discovery of algorithms using machine learning offers the prospect of reaching beyond human intuition and outperforming the current best human designed algorithms. Howeve automating the algorithm discovery procedure is intricate, as the space of possible					
Tolemic autory choose fasters from		algorithms is enormous. Here we report a deep reinforcement learning approach based on AlphaZero' for discovering efficient and provably correct algorithms forth multiplication of arbitrary matrices. Our agent, AlphaTensor, is trained to play a single-player game where the objective is finding tensor decompositions within a finite factor space. AlphaTensor discovered algorithms that outperform the state- of-the art complexity for many matrix sizes. Particularly relevant is the case of 4 × 4 matrices in a finite field, where AlphaTensor's algorithm improves on strassen's two- level algorithm for the first time, to our knowledge, since its discovery 50 years ago ² . We further showcase the flexibility of AlphaTensor through different use cases: algorithms with state-of-the-art complexity for structured matrix multiplication an improved neurical efficiency box ominizing matrix multiplication.					
	Magnetic cont through deep https://doi.org/10.1038/s41586-021-0430 Received: 14 July 2021 Accepted: ID accember 2021 Published online: 16 February 2022 Open access Check for updates	Magnetic control of tokamak p https://doi.org/10.1038/s41586-021-0430 Received: 1Decomber 2021 Published online: 10 February 2022 Open access Check for updates					

We focus on the fundamental task of matrix multiplication, and use been discovered by attacking this tensor decomposition problem using deep reinforcement learning (DRL) to search for provably correct and human search^{23,54}, continuous optimization^{17,19} and combinatorial



orization and image captioning. We believe this approach has the potential to be widely useful for better aligning models with a diverse range of

computer vision tasks.

(b) Optimize PO: $43.1 \rightarrow 46.1$, removes many incoherent predic-

tions, especially for small-scale objects.

6

Main issue with standard supervised learning: What we optimize ≠ what we care about

What do we care about?

"Good results"

Bold Metrics*

What is a metric*

$$m(x, y) \to \mathbb{R}$$
$$r(s, a) \to \mathbb{R}$$

What is a reward?

Metrics* and rewards are the same thing

Reinforcement Learning (RL) allows us to optimize what we care about

*metric in the colloquial sense. I am talking about measures

- π = neural network
- R = metric / reward func.
- S = Dataset
- s = state / image
- a = action / label

Example: Classification

1. Write down performance measure.

$$J(\pi) = \frac{1}{|S|} \sum_{(s,a) \in \mathcal{S}} R(s,a) \pi(a|s)$$

2. Compute the gradient that maximizes performance: "Policy gradient"

$$\nabla_{\pi} J(\pi) = \frac{1}{|\mathcal{S}|} \sum_{(s,a) \in \mathcal{S}} R(s,a) \nabla_{\pi} \pi(a|s)$$



J = Performance

3. Use accuracy as reward

- π = neural network
- R = metric / reward func.
- S = Dataset
- s = state / image
- a = action / label
- a* = correct class /action

$$\nabla_{\pi} J(\pi) = \frac{1}{|\mathcal{S}|} \sum_{(s,a) \in \mathcal{S}} acc(s,a) \nabla_{\pi} \pi(a|s)$$

4. Simplify equation

$$\nabla_{\pi} J(\pi) = \frac{1}{|\mathcal{S}|} \sum_{(s,a^{\star}) \in \mathcal{S}} \nabla_{\pi} \pi(a^{\star}|s)$$

5. Optimize log probability instead (doesn't change global optimum)

$$L_{\pi} = -\sum_{\substack{(s,a^{\star})\in\mathcal{S}\\ \text{Cross Entropy}}} \log \pi(a^{\star}|s) = -\log \prod_{\substack{(s,a^{\star})\in\mathcal{S}\\ \text{Negative Log-Likelihood}}} \pi(a^{\star}|s)$$

Cross-entropy is the policy gradient that maximizes accuracy!

EVERWROTEACLASSIFIERP



So how did we apply this idea in autonomous driving?

Take your driving metric and make it a reward

(3) **Driving Score (DS)**: weighted average of the route completion with infraction multiplier P_i

$$DS = \frac{1}{N} \sum_{i}^{N} R_i P_i \tag{8}$$



Also happens to scale much better than traditional reward shaping

First time we beat Supervised Learning with RL in planning



Method	Туре	Non-Reactive			Reactive			Time
		CLS ↑	Col. ↑	RC ↑	CLS ↑	Col. ↑	RC ↑	$1 \text{ me } \downarrow$
Log Replay (LQR)	Human	93.5	98.8	99.0	80.3	85.6	99.0	-
PlanTF [8]	IL	84.6	94.2	90.7	76.1	95.2	77.2	107
Diff. Planner [32]	IL	89.6	95.9	94.2	82.7	93.1	85.9	138
CaRL (Ours)	RL	91.3	97.4	94.4	90.6	97.1	91.3	14

 Table 6: Performance on Val14 (nuPlan)

Do you have an unusual metric in your problem? Maybe you should optimize it with RL.