

Inteligencja obliczeniowa

Raport 6

Temat: Implementacja algorytmów ewolucyjnych na GPU z użyciem TensorFlow

TL;DR; Dopracowałem implementację algorytmu EMAS, uwzględniając większość uprzednio uzgodnionych wymagań, co do implementacji - poza wyspami. Aktualna wersja kodu kompiluje się przez XLA, umożliwiając uruchamianie na TPU oraz przyspieszając wykonanie na pozostałych platformach obliczeniowych. Zostały dodane nowe operacje krzyżowania oraz zliczanie operacji krzyżowania. Co nieco zaskakujące, uruchomienie na TPU nie okazało się dużo szybsze od CPU.

Colab: [Final EMAS colab notebook](#)

<https://colab.research.google.com/drive/1YcMlwQgjinuP2Au-akMHnJE4dT2fqZAFE?usp=sharing>

Q&A

Co zostało zaimplementowane?

Napisana poprzednio podstawowa wersja algorytmu EMAS (Evolutionary Multi-Agent System) została dopracowana. Uzgodnione uprzednio zmiany dotyczyły następujących elementów:

- dostosować kod tak, aby był kompilowalny przez XLA,
- dodać do implementacji wyspy oraz mechanizm migracji agentów między nimi,
- zrobić zliczanie operacji krzyżowania,
- dodać nowe operacje krzyżowania.

Punkty zaznaczone kolorem zielonym zostały poprawnie wykonane, natomiast ten na czerwono - nie.

Co zostało zrobione w ramach umożliwienia kompilacji przez XLA?

Głównym problemem była zmienna długość tensorów w kolejnych iteracjach pętli while. Aby poradzić sobie z nim zostało zastosowane wykorzystanie pętli while w formie explicite oraz dopełnianie problematycznych tensorów do prespecyfikowanej długości wyznaczonej przez maksymalną dozwoloną wielkość populacji. Ponadto operacja shuffle wykorzystana do mieszania populacji została zastąpiona prostszą wersją opartą o resztę z dzielenia.

```
perm = (tf.cast(tf.range(current_pop_size), dtype=tf.int32) * 24179) % current_pop_size  
energy = tf.gather(energy, perm)
```

Jak zostały dodane wyspy oraz mechanizm migracji?

Nie zostały dodane. Okazało się, że dostosowanie kodu pod XLA zajęło mi więcej czasu niż się spodziewałem, sprawiając, że byłem już zbyt zdemotywowany, żeby na poważnie przysiąść do tego punktu.

Jakie operacje krzyżowania zostały dodane?

```
19 @tf.function()
20 def crossing(a, b, n):
21     k = tf.reduce_sum(tf.ones_like(a[:, 0]), dtype=tf.int32)
22
23     # Random bubble (hyper-sphere) between points
24     reposition = tf.random.uniform([k, n], minval=-1, maxval=1)
25     reposition = reposition / tf.norm(reposition, axis=1, keepdims=True)
26     new_position_1 = ((a + b) / 2.0) + (tf.maximum((tf.norm(a - b, axis=1, keepdims=True) / 2.0), 1e-3) * reposition)
27
28     # Random exchange of genomes
29     new_position_2 = tf.where(
30         tf.random.uniform(shape=[k, 1]) < tf.random.uniform([1], 0.05, 0.95), a, b)
31
32     # Choose one crossing or another
33     new_position = tf.where(
34         tf.random.uniform(shape=[k, 1]) < tf.random.uniform([1], 0.2, 0.5), new_position_1, new_position_2)
35
36     # Apply mutation, because why not
37     # Cauchy distribution derived from uniform distribution
38     mutation = 0.05 * tf.math.tan(tf.random.uniform(minval=-1.0 * math.pi / 2 + 1e-6, maxval=math.pi / 2 - 1e-6, shape=[k, n]))
39     new_position_m = tf.where(
40         tf.random.uniform(shape=[k, 1]) < tf.random.uniform([1], 0.05, 0.4), new_position + mutation, new_position)
41
42     # Choose mutated or not
43     new_position = tf.where(
44         tf.random.uniform(shape=[k, 1]) < 0.4, new_position_m, new_position)
45
46     new_position = tf.clip_by_value(new_position, clip_value_min=-5.12, clip_value_max=5.12)
47
48     return new_position
```

1. Losowanie punktów z wnętrza hiperkuli będącej pomiędzy obiema agentami
2. Wybór elementu genomu losowo z jednego bądź drugiego agenta
3. Mutacja losowana z rozkładu Cauchy'ego

Każda z powyższych operacji ma oczywiście pewną, raczej rozsądną (choć to pewnie hiper-parametr, którym warto by się zająć), szansę zajścia.

Jaką funkcję optymalizowałem?

Funkcja Rastrigina dla n-wymiarów. Testowałem w wersji lekkiej 10-D, ale dłuższej (250,000 epok), oraz cięższej 10000-D, ale krótszej (100 epok).

Jakie hiper-parametry symulacji przyjąłem?

Liczba agentów: 100

Początkowa energia: 5 dla każdego agenta

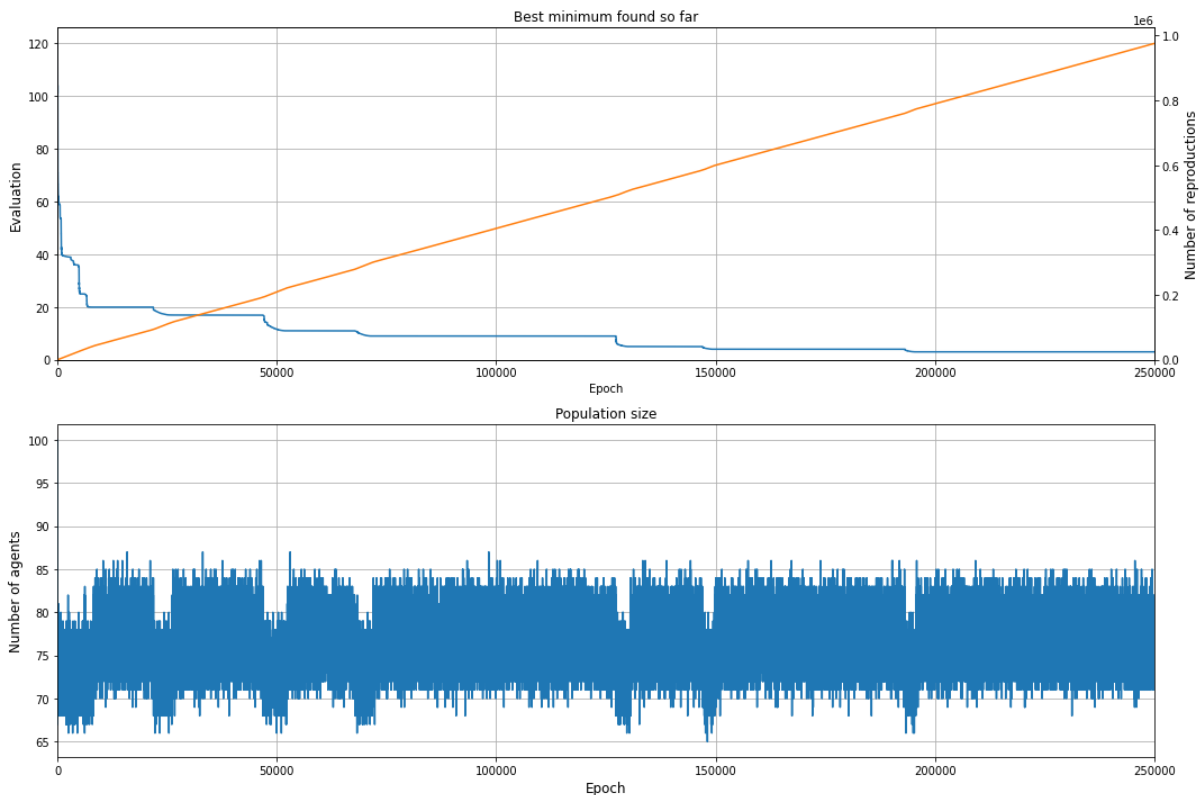
Zmiana energii po walce: 2 (przeegrany -2, wygrany +2)

Energia wymagana do reprodukcji: 10

Koszt reprodukcji: 2.5 (obu agentów -2.5, nowy agent startuje z 5)

Czy implementacja okazała się poprawna i zbiegała?

Wygląda na to, że tak. Poniżej przedstawione są statystyki uruchomienia dla 100 agentów w populacji, dla 10-wymiarowej funkcji Rastrigina w obrębie 250000 epok.



Pierwszy wykres przedstawia zmianę najlepszego rozwiązania w czasie oraz liczbę krzyżowań. Ewaluacje są wykonywane tylko w obrębie areny walk.

Drugi wykres przedstawia zmianę rozmiaru populacji w czasie. Jak widać, dla zadanych hiper-parametrów populacja bardzo szybko osiąga stabilny rozmiar około 70-85 osobników, oraz widać momenty, w których w puli genetycznej pojawiły się interesujące mutacje.

Dla pewności, suma energii w systemie jest kontrolowana przy każdym uruchomieniu na sam koniec procesu ewolucyjnego.

```
assert tf.abs(tf.reduce_sum(energy) - k * 5) < 1e-4, "Energy not preserved"
```

Jak wygląda kwestia czasu obliczeń?

Rozmiar problemu: niewielki, ale długi

- liczba agentów: 100
- liczba wymiarów: 10
- liczba epok: 1000

Rozmiar problemu: duży, ale krótki

- liczba agentów: 2000
- liczba wymiarów: 10000
- liczba epok: 100

Przez “bez XLA” w poniższej tabelce rozumiemy, że zostały wykorzystane wszystkie optymalizacje poza kompilacją XLA (jit_compile).

	Niewielki, ale długi		Duży, ale krótki	
	bez XLA	XLA	bez XLA	XLA
CPU	14.9583s	12.1042s	64.3765s	60.5606s
GPU - Standard	43.2964s	38.9117s	0.5612s	0.4856s
GPU - Premium	52.3314s	50.5861s	0.1559s	0.1379s
TPU	niewykonalne	11.5694s	niewykonalne	58.9382s

Ogromnym zaskoczeniem jest wynik dla TPU, który jest bardzo zbliżony do czasów dla zwykłego CPU. Ciężko określić czy obliczenia były faktycznie wykonywane na TPU (powinny), czy też w tym wypadku z jakiegoś powodu mimo wszystko wykonało się na CPU. Jeśli faktycznie wykonało się to na TPU, to możliwe, że kompilator postanowił zlinearyzować kod na wzór tego dla CPU, kompletnie nie wykorzystując potencjału do akceleracji - notabene czasy były powtarzalnie lepsze (choć tylko odrobinę) od CPU, ale nie mam tego nigdzie formalnie udowodnione dla większej próby, więc to mógł być też przypadek.

W ogólności natomiast widać, że po kompilacji XLA kod na wszystkich maszynach jest faktycznie nieznacznie szybszy od kodu, który miał wszystkie optymalizacje poza tą jedną.

Wnioski z projektu

Pisanie kodu kompilowalnego do XLA w TensorFlow, o ile raczej łatwiejsze i bardziej wysokopoziomowe niż pisanie w czystej Cudzie, nadal pozostaje sporym wyzwaniem, gdyż pojawia się dużo nieoczywistych problemów:

- niektóre operacje nie są wspierane, trzeba mieć świadomość wielu ograniczeń,
- debugowanie jakiegokolwiek dłuższego kodu jest wyzwaniem samym w sobie:
 - wiele błędów ujawnia się dopiero, gdy dodana zostanie flaga uruchamiająca kompilację kodu, bez której może działać idealnie,
 - na stosie błędów wywołania są do funkcji ze zbudowanego grafu obliczeniowego, a nie oryginalnego kodu.

Mimo to, z odpowiednią dozą samozaparcia i kreatywności możliwe jest wykonanie dosyć skomplikowanych operacji - w ramach projektu były to: wcześniej algorytm mrówkowy, oraz teraz algorytm EMAS.

Pod względem szybkości działania takiego kodu, trzeba zaznaczyć, że testy wykonywane w trakcie projektu pokazywały przyspieszenia na poziomie nawet dwóch rzędów wielkości, względem czystego pythona wywołującego operacje interaktywnie operacje tensorflow. Widoczne jest również przyspieszenie względem funkcji, dla których wyznaczony zostaje graf obliczeń (flaga tf.function).