

MACHINE LEARNING

Practical work 04 Genetic algorithm



HEIG-VD
Burgener François

1 Introduction

Le problème du voyageur de commerce est assez connu. Il demande de trouver un parcours entre chaque ville le plus court possible en revenant au point d'origine. Le parcours ne doit jamais passer plus d'une fois sur une ville mise à part la ville de départ à la fin du parcours. Ce problème est un problème difficile. Ce problème est un NP-complet. Il nous est donc pas toujours possible de trouver avec certitude la solution optimale. En théorie il serait possible de trouver la solution optimale car nous pouvons avoir tous les chemins possibles mais cela nous prendra beaucoup trop de temps pour certains problèmes plus complets.

Dans notre cas nous utilisons ce problème sur 14 villes en Birmanie via des coordonnées latitude, longitude. Pour le calcul de la distance nous avons utilisé la formule de Vincenty. Cette formule nous permet de calculer la distance entre deux points à la surface d'une sphère. Ces formules utilisent l'hypothèse que la terre est une sphère aplatie aux pôles ce qui permet d'obtenir des résultats plus précis qu'avec la formule de la distance du grand cercle.

2 Solution

Pour notre solution nous avons utilisé l'algorithme génétique afin d'éviter de devoir coder la solution du problème. L'algorithme nous permet de trouver la solution sans devoir coder le problème en lui-même. Avec cet algorithme, nous pouvons définir quel type de solution nous souhaitons obtenir afin qu'il trouve par lui-même la solution optimale. Le type de solution que nous avons définie est une liste de parcours de chaque ville. L'algorithme va prendre les listes et prendre la meilleure, qui pour nous est le chemin le plus court. Nous utilisons donc une fonction d'évaluation qui va calculer la distance du parcours en entier afin de pouvoir faire cette minimisation.

3 Provide the better route you found and the shortest path in kilometers. Is it the optimal shortest path ?

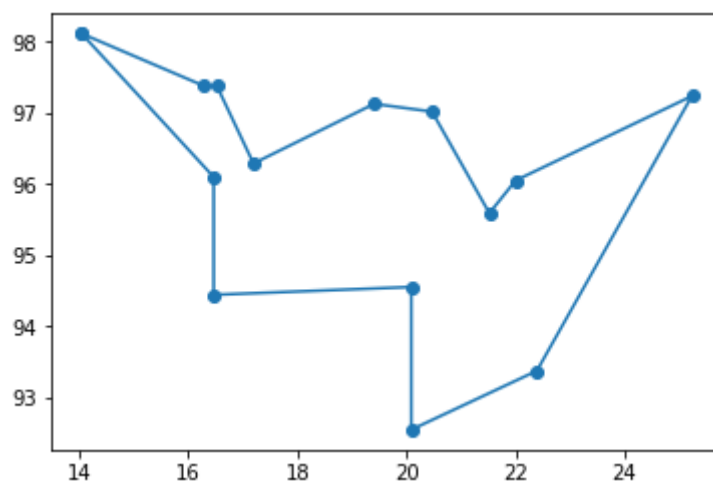
Meilleur parcours : [9, 0, 1, 13, 2, 3, 4, 5, 11, 6, 12, 7, 10, 8].

Nous pourrions avoir d'autres parcours similaires qui seraient juste une rotation ou une inversion de notre solution. Le parcours sera identique et donc la distance sera la même.

Distance du meilleur parcours : 3346.76 km

Notre parcours est le plus optimal. Nous arrivons à trouver le plus court chemin en à peine 80 générations. J'ai fait des tests en allant jusqu'à 10'000 générations et les résultats étaient toujours les mêmes. Par contre, il aurait été possible que notre solution ne soit pas optimale car dans certains cas, il n'est pas possible de trouver une solution optimale (je parle dans un autre exemple que les 14 villes données).

Résultat du chemin le plus court :



4 Fonction Fitness

Notre fonction fitness prend en paramètre un chromosome. Ce chromosome est une liste des villes parcourues. Ensuite, nous calculons, en kilomètre, la distance du parcours de toutes ses villes en retournant à la ville d'origine. Ensuite, nous retournons cette distance en tant que score. Ce score permettra à l'algorithme génétique de savoir si nous avons le meilleur parcours ou non. Il gardera le meilleur parcours en fonction de si on maximise ou minimise le résultat du score. Dans notre cas, nous voulons minimiser la distance du parcours.

5 Encoding solution

Pour coder la solution j'ai décidé de prendre une liste d'entier représentant le parcours effectuer dans chaque ville. Cette liste a des chiffres entre 0 et 13.

Voici un exemple de chromosome possible : [9, 0, 1, 13, 2, 3, 4, 5, 11, 6, 12, 7, 10, 8]

Ces numéros nous permettront de récupérer les coordonnées (latitude, longitude) de chacune des villes car nous avons un tableau de coordonnées.

Tableau des coordonnées : [(16.47, 96.1), (16.47, 94.44), (20.09, 92.54), (22.39, 93.37), (25.23, 97.24), (22.0, 96.05), (20.47, 97.02), (17.2, 96.29), (16.3, 97.38), (14.05, 98.12), (16.53, 97.38), (21.52, 95.59), (19.41, 97.13), (20.09, 94.55)]

Comme ça dès que l'on récupère une valeur dans notre chromosome nous pouvons prendre ce numéro (qui est le numéro de la ville) et l'insérer en tant qu'index dans le tableau des coordonnées afin de récupérer ses coordonnées. Ensuite il nous suffit d'utiliser ses coordonnées dans notre fonction de calcul de distance pour connaître la distance entre deux villes.

Pour initialiser notre chromosome, j'utilise la méthode proposée dans l'exemple TSP de pyevolve. Cette méthode nous crée une liste de la taille de notre génome allant de 0 à 13. Ensuite nous mélangeons cette liste afin d'avoir une liste aléatoire.

Pour ce qui est de la mutation j'utilise la méthode de swap qui va échanger aléatoirement des éléments d'un génome. Pour ce qui est du crossover je ne l'utilise pas même si je l'ai initialisé dans la solution. En ayant fait des tests avec j'ai pu remarquer que de ne pas l'utiliser nous donne une meilleure performance. Nous trouvons la solution optimale moins de générations qu'avec le crossover rate.

Pour cela j'ai utilisé le notebook **LaboGA.ipynb**.

6 Configuration of the GA

Pour la solution finale j'utilise comme paramètre :

Generation : 500

Mutation : 0.1

CrossoverRate : 0.0

Population : 10

Selector : GrouletteWheel

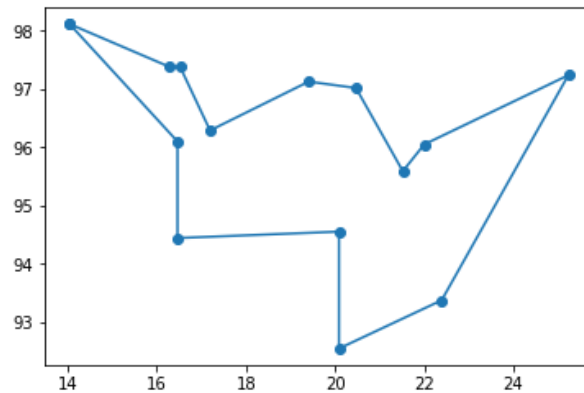
Elitism : True

Comme je l'ai décrit plus haut, j'utilise la méthode swap pour les mutations. En ce qui concerne le crossover je ne l'utilise pas et je le mets à 0.0 ce qui signifie qu'il est set à 0%.

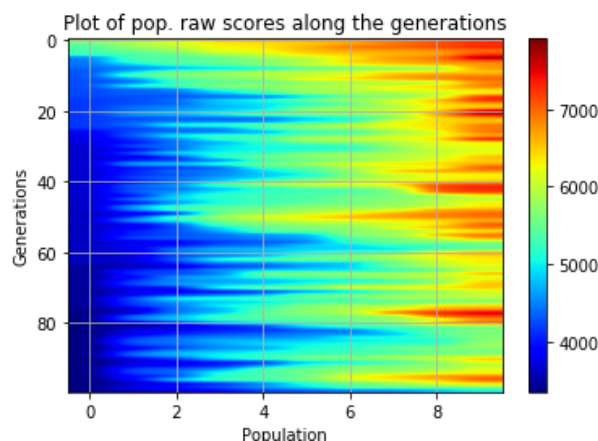
Au tout début de mes tests j'utilisais une population élevée afin de faire mes tests. Il me fallait donc un grand nombre de générations afin d'arriver à la solution optimale. J'ai donc préféré avoir une plus petite population afin de récupérer les meilleurs résultats plus rapidement. Ensuite comme je n'utilisais pas le crossoverRate je n'avais que la mutation qui changeait. J'ai utilisé différents pourcentages de mutation afin d'arriver rapidement à la solution optimale. En ce qui concerne la génération j'ai préféré le laisser à 500 afin d'éviter de trouver la solution optimale en beaucoup de générations.

7 Provide relevant plots of your experiments and explanations

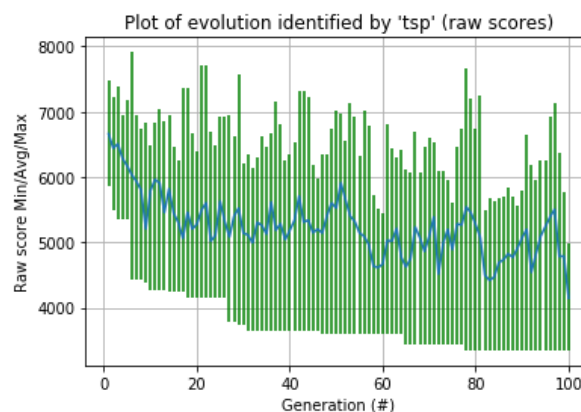
Ici nous pouvons voir le parcours le plus court que nous avons obtenu. Chaque solution que nous trouvons aura le même résultat que le graphe si dessous. La seule différence est le sommet de départ et la direction dans le quel on parcourt les ville. Soit dans le sens des aiguilles d'une montre ou l'inverse. Cela nous donne donc 28 solutions de parcours (14 dans un sens et 14 dans l'autre).



Ce graphe nous montre une carte de la distribution des scores brute de la population entre les générations. Ici nous pouvons voir qu'à partir de la génération 80 nous avons un grand nombre de la population qui est dans le bleu ce qui est pour nous une bonne chose car nous voulons que le score soit le plus bas possible c'est-à-dire dans le bleu.



Ce graphe nous permet de voir l'évolution du score entre la génération. Ici nous allons regarder le minimum est on peut voir qu'à partir de 80 le score ne change pas est qu'il reste au même stade. Cela confirme les résultats obtenus par l'algorithme génétique.



8 Conclusion

Pour conclure, nous pouvons remarquer que le problème TSP, avec le set de données qu'il ou était fournie, était un problème simple pour trouver la solution optimale. Avec la bonne configuration de notre algorithme générique, que se soit sur l'initialisation, les mutations ou alors le crossover rate, nous pouvons trouver très rapidement, et ça en peu de génération, la solution optimale. Au tout début j'avais de la peine à comprendre comment la GA pouvait trouver par lui-même une des solutions optimales. En fait c'est grâce à notre fonction d'évaluation (fitness) qui permettait à l'algorithme génétique de dire si le chromosome choisi doit rester ou non dans la population. Si ce chromosome avait un meilleur score que les autres de la population, il était gardé et le plus mauvais chromosome était supprimé de la population. Cela permet de garder que les meilleurs individus dans la population.