

Optimization des Hyperparamètres appliquée au Fine Tuning de LLM

Basé sur l'article : *Bayesian and Partition-Based Optimization for Hyperparameter Optimization of LLM Fine-Tuning*

Nathan Davouse

Sommaire

1. Introduction

2. Design et Implémentation

3. Résultats et Analysis

4. Conclusion

Large Language Models

Point clés

- Etat de l'art pour le traitement de langage naturel.
- Réseaux de Neurones avec une architecture basé sur le transformer^a
- Taille : entre 1 et 405 Milliards de neurones

^aVaswani et al, Attention is all you need, 2017

Auto-attention

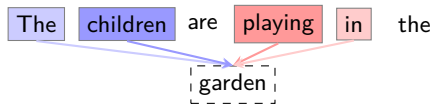


Figure: Illustration du mécanisme d'auto-attention

L'auto-attention est la clé du LLM, en permettant de comprendre le contexte

Fine Tuning

- second phase de l'entraînement (prendre que la partie de droite)

PEFT

Utilisation de la méthode LoRA, qui permet de réduire les couts d'entraînement. (cf annexe 2)

Optimisation des Hyperparamètres (OHP)

Hyperparamètres

Paramètres qui ne sont pas entraînés par le modèle
(learning rate, dropout ...)

Objectifs

- ▶ Meilleure performance qu'en manuel
- ▶ Retirer le besoin d'expertise

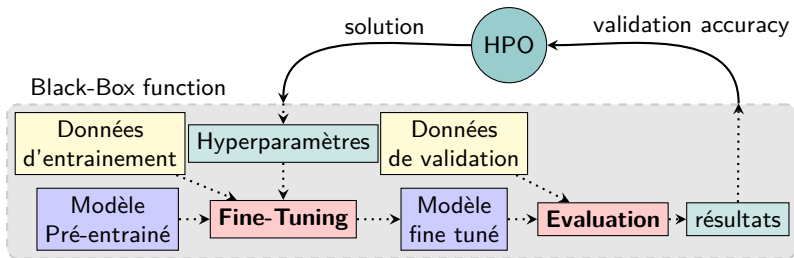


Figure: Fonctionnement général de l'optimisation des hyperparamètres

Sommaire

1. Introduction

2. Design et Implémentation

3. Résultats et Analysis

4. Conclusion

Search Space

Search Strategy : B0

Search Strategy : S00

Search Strategy : BaMSOO

Performance Estimation Strategy

Implémentation

Sommaire

1. Introduction

2. Design et Implémentation

3. Résultats et Analysis

4. Conclusion

Expérimentation

LHS : Résultats

Résultats des 3 algorithmes

Analyse

Prospectives

Sommaire

1. Introduction

2. Design et Implémentation

3. Résultats et Analysis

4. Conclusion

Conclusion

Une conclusion

Merci.

Annexes 1 : Architecture d'un LLM

MHA, Transformers

Annexes 2 : Low Rank Adaptation (LoRA)