

# LoRA

---

**Hyunbeen Kim**

been000904@gmail.com

**NLP Team**

2024/05/28



# Contents

---

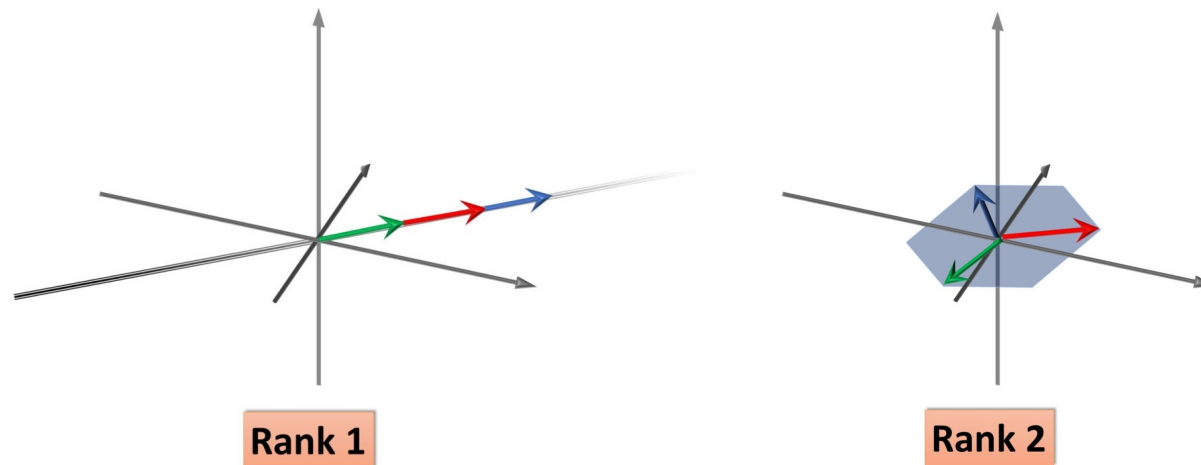
- Prerequisite
- Simple Summary
- Motivation
- Method
- Advantage

# Prerequisite

---

## Rank

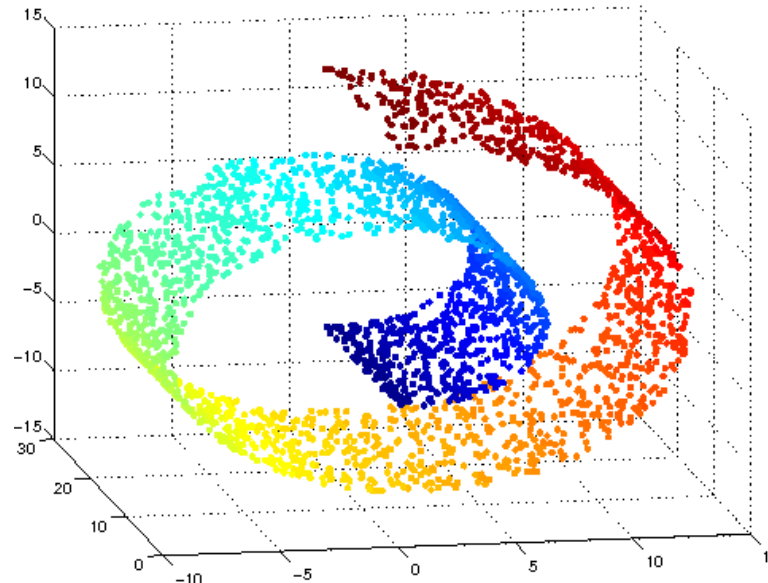
In linear algebra, the rank of a matrix  $A$  is the dimension of the vector space generated (or spanned) by its columns.



# Prerequisite

## Intuition to Model Rank

데이터가 표현되기 위해 필요한 차원만큼 모델이 구성할 수 있어야 한다.  
필요한 차원 이상으로 모델이 표현력을 가지게 되면 오버피팅,  
필요한 차원 이하로 모델이 표현력을 가지게 되면 언더피팅



# Simple Summary

---

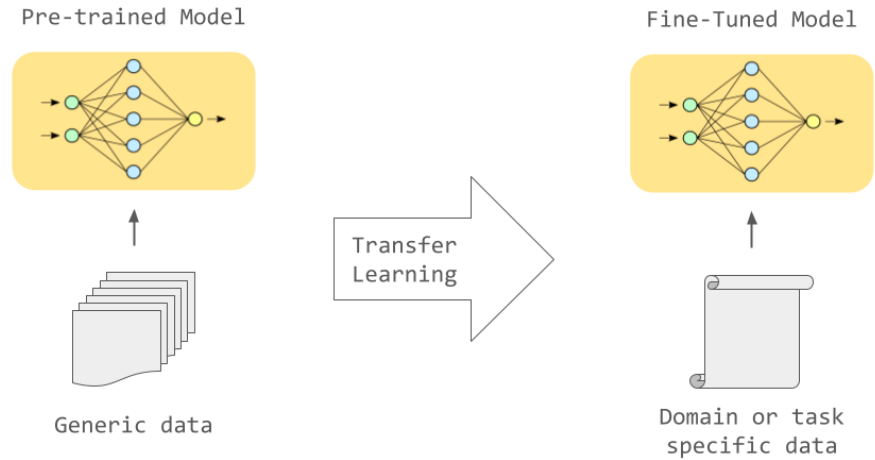
일반적인 도메인에 대한 이해를 가지고 있을 때, 그보다 작은 도메인에 대해 새로운 이해를 하기 위해 필요한 파라미터의 랭크는 절대 크지 않다. 즉, 낮은 랭크를 가지는 추가 파라미터만 튜닝함으로써 작은 도메인에 대한 이해를 가지도록 만들 수 있다.

이는 LLM을 포함한 모든 Neural Network에 적용할 수 있다.

# Motivation

---

Language Model에서 주류는 일반적인 도메인에서 커다란 모델을 학습하고  
이를 필요한 Task에 맞춰 재학습하는 것, 보통 파인튜닝을 이용하여 진행



# Motivation

---

파인튜닝으로 접근하는 경우, 모델 학습 비용도 굉장히 많이 들고, 간단한 문제를 풀 때나 복잡한 문제를 풀 때나 동일한 모델 아키텍처를 가져가게 되는 단점 존재 (간단한 문제를 커다란 모델로 푸는 비효율)

# Motivation

---

기존의 논문을 통해 생각해보았을 때 새로운 지식 습득하는 과정(Fine Tuning)을 모델링하기 위해서 꼭 높은 랭크를 가지는 파라미터가 필요한 것은 아니고 낮은 랭크 (low intrinsic rank)로도 충분히 모델링할 수 있음



# Motivation

---

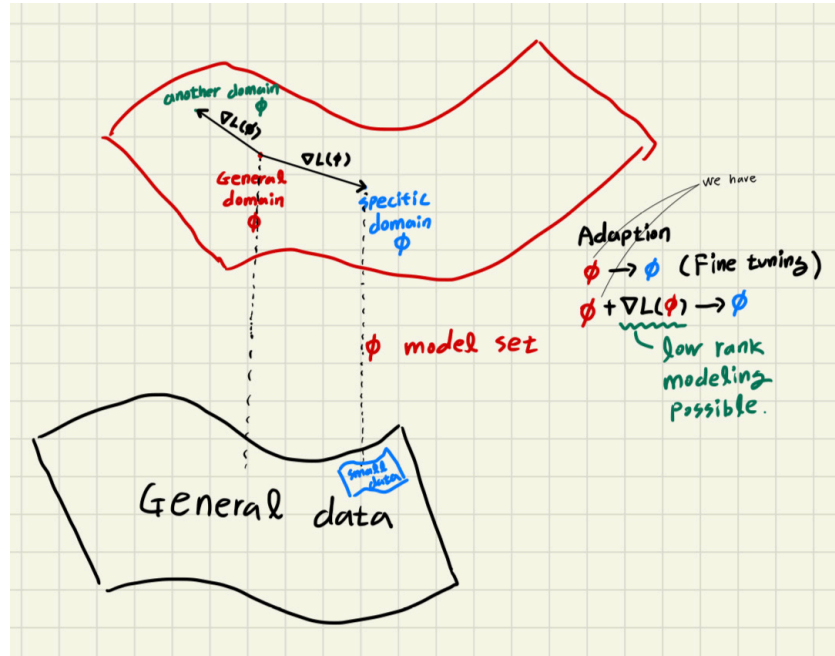
즉 MLE를 할 때 전체 파라미터를 추정하는게 아니라, 변화량만 추정하면 어떨까?  
전체 파라미터는 가만히 놔두고 변화량을 모델링하는 파라미터만 학습

$$\max_{\Phi} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log (P_{\Phi}(y_t|x, y_{<t}))$$

$$\max_{\Theta} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log (p_{\Phi_0 + \Delta\Phi(\Theta)}(y_t|x, y_{<t}))$$

# Motivation

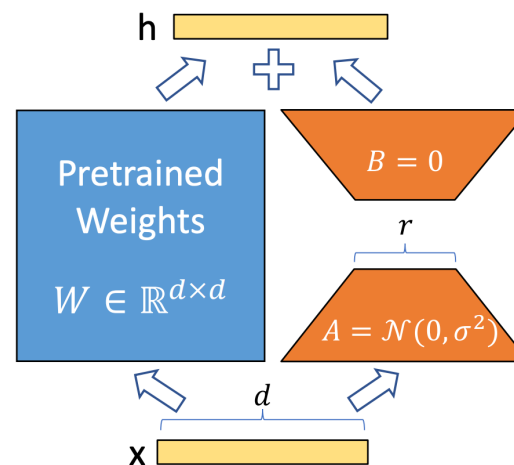
즉 MLE를 할 때 전체 파라미터를 추정하는게 아니라, 변화량만 추정하면 어떨까?  
전체 파라미터는 가만히 놔두고 변화량을 모델링하는 파라미터만 학습



# Method

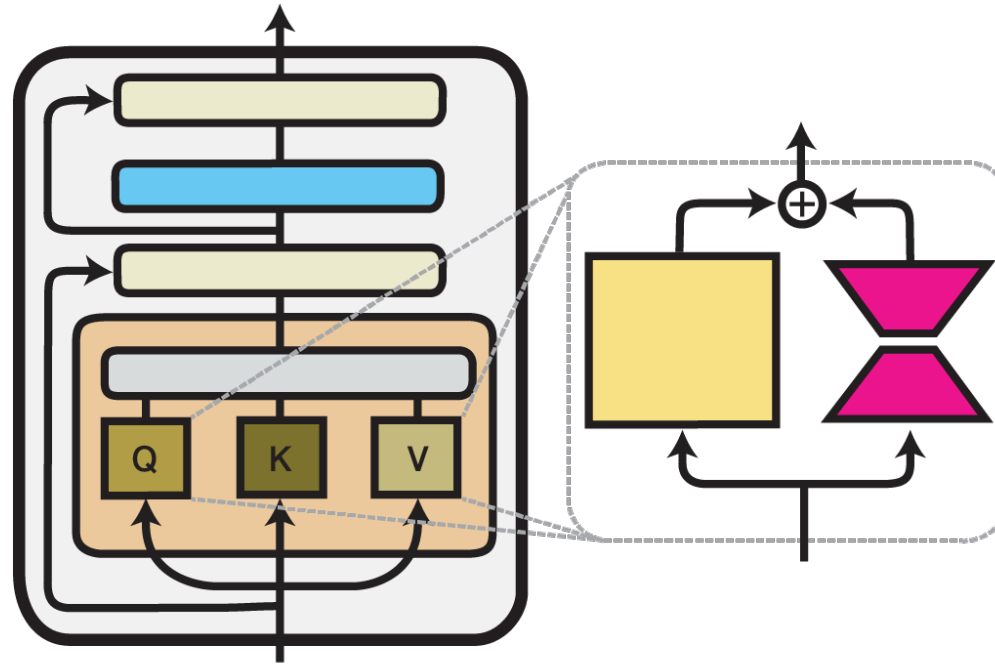
어떻게 변화량을 작은 Rank로 근사할 수 있을까?

$$h = W_0x + \Delta Wx = W_0x + BAx$$



# Method

그럼 Low Rank로 변화량 모델링하는걸 어떻게 쓸까?



# Advantage

---

1. Training Cost를 압도적으로 줄임, 가령 GPT-3 175B의 0.01%에 해당하는 파라미터만 튜닝하여 Downstream Task에 적용 성공
2. Inference 할 때 Latency를 전혀 늘리지 않고 사용할 수 있고, 프로덕션에 배포할 때 적은 가중치만 바뀌가면서 다양한 Task를 수행할 수 있으니 장점을 가짐
3. 원본 가중치를 훼손하지 않으니 Prior Knowledge를 많이 가져갈 수 있음



TRAIN AND TEST