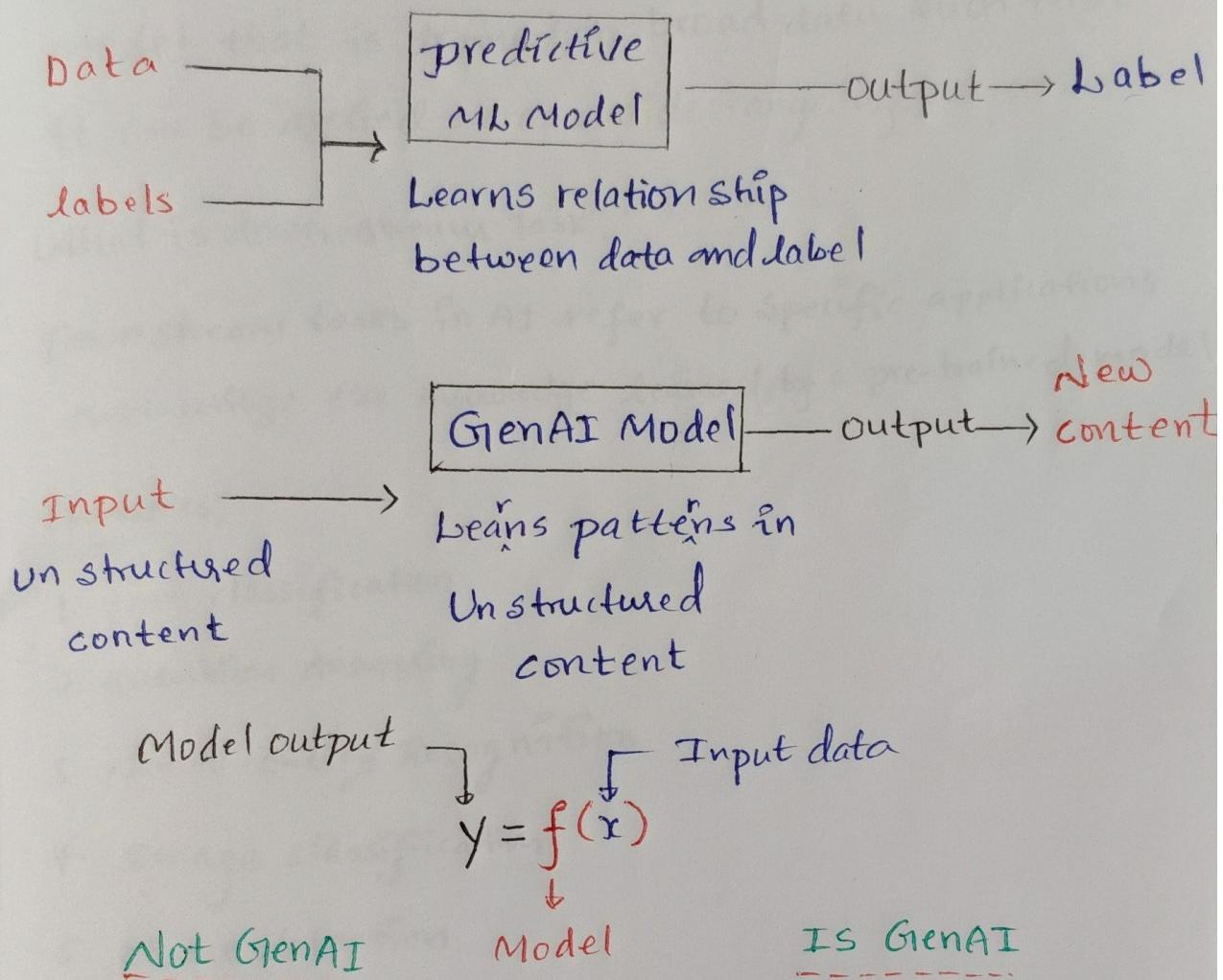


Generative AI

Date: 30-06-2024

Def: GenAI is a type of Artificial Intelligence that creates new content based on what it has learned from existing content.

How to Understand GenAI?



When y is a:

- Number
- Discrete
- Class
- probability

When y is:

- Natural language
- Image
- Audio
- Video

What is fine-tuning?

Making small adjustments to foundation model in-order to achieve the best performance in a down-stream task.

What is foundation model:

Is also known as large AI model, is a ML or DL model that is trained on broad data such that it can be applied across wide range of use-cases.

What is down-stream Task:

Downstream tasks in AI refer to specific applications that utilize the knowledge learned by a pre-trained model

examples:

1. Text classification
2. Question Answering
3. Named Entity Recognition
4. Image classification
5. Object detection
6. Medical Image Analysis
7. Fraud detection
8. Algorithmic Tracking

- If we update all the model parameters (Re-train) during finetuning then it is called **full finetuning**

Parameter Efficient Fine-Tuning (PEFT)

Why We Need PEFT:-

The major downside of finetuning is that the new model contains as many parameters as in a original model.

- As larger models contains as many trained every few months, this changes from a mere "inconvenience" for GPT-2 to a critical deployment challenge for GPT-3 with 175 billion parameters

Low Rank Adaptation (LORA):

- LORA is one of the most used fine-tuning technique in the umbrella of PEFT

AREN'T EXISTING SOLUTIONS GOOD ENOUGH?

- There are two prominent strategies when it comes to efficient adaptations:
 1. Adding Adapter layers.

2. Optimizing some forms of the input layer activations

1. Adapter layers introduce **Inference Latency**:

There are many variants of adapters. Adapter layers are those which we fine tune / train for the downstream task.

- There are no direct ways to bypass extra compute in adapter layers. Since we increased no. of parameters by adding the extra layers inference latency is increased.

2. Directly optimizing the prompt is **Hard**:

Observed that prefix tuning is difficult to optimize and that its performance changes non-monotonically in trainable parameters.

prefix tuning: It is nothing but adding a prefix before the prompt.
before

Eg: prompt: "Generate a smart phone product description".

prompt_prefix: [positive tone] generate a product description for a smartphone".

OUR METHOD: (LORA)

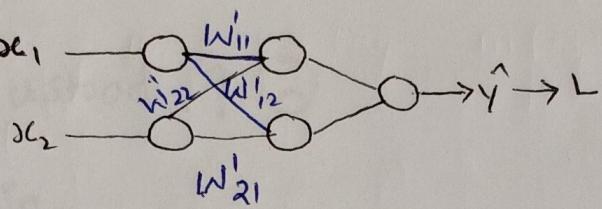
First, what is Rank ?

The no. of independent rows or columns in a matrix.

- Do, explore the examples 😊

* We know that weights in a specific layer of neural networks can be written in matrix form.

Example:



$$W^1 = \begin{bmatrix} w_{11}^1 & w_{21}^1 \\ w_{12}^1 & w_{22}^1 \end{bmatrix}_{2 \times 2}$$

* When we are adapting to a specific task, we came to know that the pre-trained language models have a low "Intrinsic dimension" (Search for research paper, you will get it on intrinsic dimension).

Wait! Wait! Wait! What is Intrinsic dimension

Jayanth?

An objective function's **intrinsic dimension** measures the minimum number of **parameters** (weights + biases) needed to reach satisfactory solutions for respective objective (or)

The **Intrinsic dimension** represents the lowest dimensional **Subspace** in which one can optimize the original **objective function** to within a certain level of approximation error.

Not Understood ?? 😔

Let me explain,

Let's take GPT 3, which has 175 Billion parameters.

Let $\Theta^D = [w_1, w_2, \dots, w_{175\text{Billion}}]$



Consider the subspace of foundational model parameters

Let $\Theta^d \subseteq \Theta^D$ [$d \leq D$]

Basically think Θ^d is the projection of Θ^D

Let $\Theta^d = [w_1, w_2, \dots, w_{1\text{million}}]$

think Θ^d is the **basis** of Θ^D

There will be so many basis for θ^D , which one I have to select Jayanth?

Well,

The authors are saying that:

We have to select θ^D such that it gives 90% of full-finetuning accuracy.

Then they are calling this θ^D as Intrinsic dimension

- σ^2

Now,

Taking the inspiration from Intrinsic dimension the LORA authors thought there would be Intrinsic Rank as well for LLM's.

* Author's hypothesized the updates to weights also have a low "Intrinsic rank" during adaptation

Let $W_0 \in \mathbb{R}^{d \times K}$ is the pre-trained weight matrix (foundation model) of a layer

Let $\Delta W = B_{d \times r} \cdot A_{r \times K}$, where

'B' is initialized with "zeros" &

'A' is Random Gaussian Initialization"

Now instead of training W_0 (frozen) we train A and B matrices for the down-stream task.

low-Rank decomposition

Then we multiply them to get ΔW , since ΔW is same dimension as W_0 . We add them $d \times K$

$$h = x^{W_0} + \Delta W x = W_0 x + BA x$$

↑
Input (embedding · prompt)

<show figures>

Jayanta, we apply LORA on a foundation model right but you apply on a layer?

yes, when we apply LORA on a foundation model like llama. It means we are applying LORA to each layer of it!

Remember,

before we add $\Delta W x$ with $W_0 x$ we scale it with $\frac{d}{K}$.

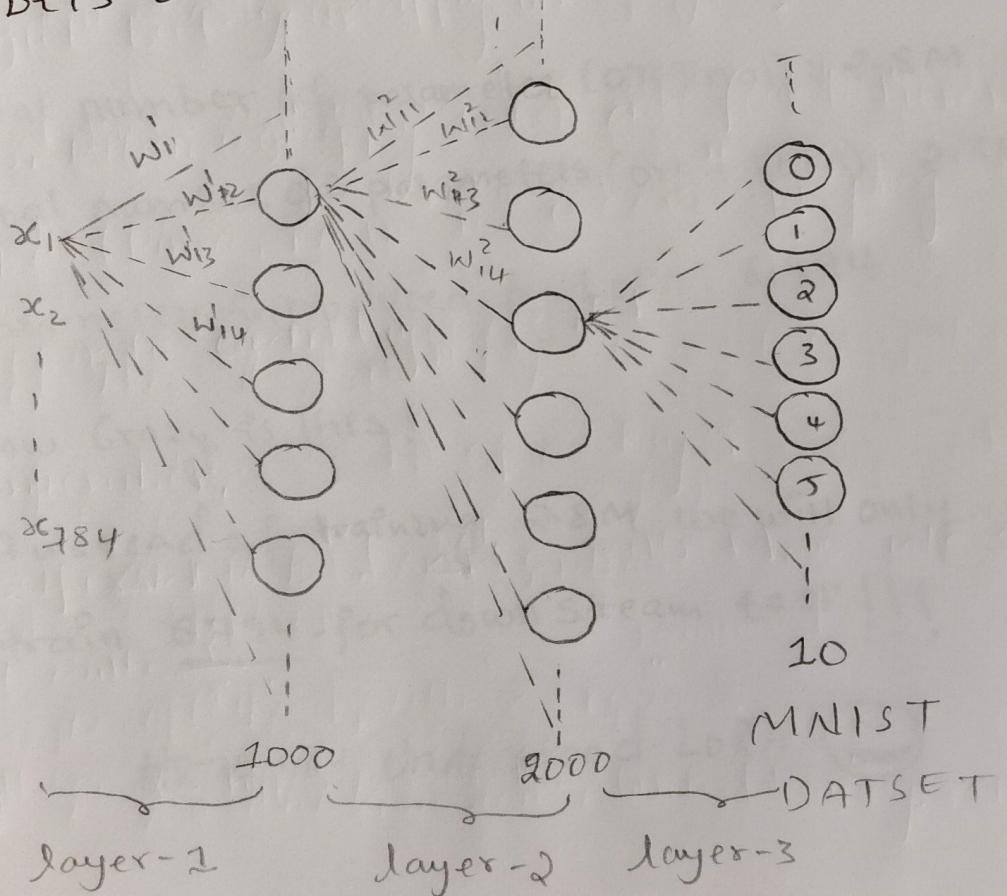
ultra. α -hyperparameter like learning rate

Why we need to scale $\Delta w \times ?$

- Controlling the influence of LORA updates:
 - A higher α/γ means the lora updates will have a stronger effect.
 - lower $\frac{\alpha}{\gamma} \Rightarrow$ keeping the model closer to the original pre-trained weights.
 - Stabilizing the Training process:

when using higher LORA ranks. Without this scaling, the LORA updates could overwhelm the original model weights and lead to **unstable training**.

Let's take an example:



Output Input

↑ ↓

Layer 1: $W_0: [1000, 784] + B_0: [1000]$

Layer 2: $W_0: [2000, 1000] + B_0: [2000]$

Layer 3: $W_0: [1000, 2000] + B_0: [10]$

Let's Apply LORA to Layer 1:

Layer 1: $W_0: [1000, 784] + B_0: [1000] +$
 $\text{loraA}: [1, 784] + \text{loraB}: [1000, 1]$

Now,
 we freeze W_0 & B_0 and train only LORA-A
 & LORA-B matrices for Layer-1

likewise for other layers

Total number of parameters (original): 2.8M

Total number of parameters (org + LORA): 2.813M

parameters introduced by LORA: 6,794

How crazy is this!

Instead of training 2.8M we will only
 train 6,794 for downstream task!!!

Hope you understood LORA 😊