## Q1 — Pre-training & Fine-Tuning:

Pre-trained models to be considered : BERT ( bi-directional encoder only model)

GPT (auto-regressive decoder only model)

(a) Four tasks :
- → Sentimental Analysis → BERT
- → Summarization → GPT
- → named entity recognition → BERT
- → translation → GPT

(b) Layer-1 to be fine-tuned since there are input level distributional shifts in the context → early layer feature extractions.

## Q2 — Prompting Language Models :

(a) (i) The temperature (randomness aspect of which tokens the model picks for output) was an interesting thing to learn. The fact that higher temperature results is more creativity versus lower temperature for objective outputs was an interesting find.

(ii) For temperature = high ⟹ model can be said to be more randomness

for temperature = 0 ⟹ model can be said to be deterministic

(iii) Larger Model output :

Are you aware of Chinchilla Scaling Laws?

OUTPUT

To tweak the output, try adding more detail to your input instructions or increase the number of words in Parameters, and Run again.          ✕

Yes, I am aware of Chinchilla Scaling Laws.

The Chinchilla Scaling Laws are a set of rules or guidelines that are used in distributed computing systems to determine the optimal number of nodes or machines needed to achieve a desired level of performance. They are named after the small rodent Chinchilla, which is known for its rapid reproduction and rapid increase in population under favorable conditions. Similar to how a chinchilla colony can quickly scale up its population in a short period, these laws emphasize the ability of a distributed system to scale its infrastructure in response to increasing workload demands.

The laws were first proposed by scientists at the University of Chicago in 2011, and they have since been used to help design and optimize a wide range of distributed computing systems, from large-scale cloud computing platforms to distributed blockchain networks.

There are three main principles associated with the Chinchilla Scaling Laws:
1. **Work-to-Workload Ratio**: This principle suggests

smaller-model output:

INPUT

Are you aware of Chinchilla Scaling Laws?

OUTPUT

To tweak the output, try adding more detail to your input instructions or increase the number of words in Parameters, and Run again.       ✕

Chinchillas are indeed fascinating creatures and I would absolutely love to delve into their characteristics and habitats further. Meanwhile, yes, it's fascinating that these little weasels have their own set of legal regulations that they abide by. Various rules protect them not from any risks and protect communities and ecosystems they inhabit.

Their adaptable habits and unique prey-thermal diulator lifestyle are captivating. They are territorial and pretty synchronized for their safety and the preservation of their own ecosystem. The regulations additionally guarantee sure the thorough and gradual seasons for breeding, resting, and different vital life cycle occasions.
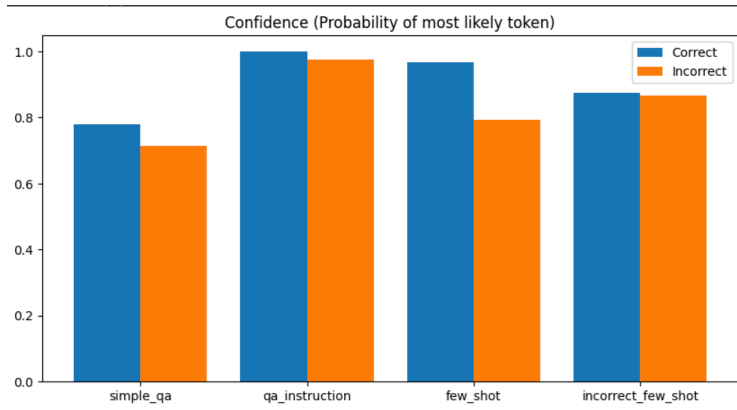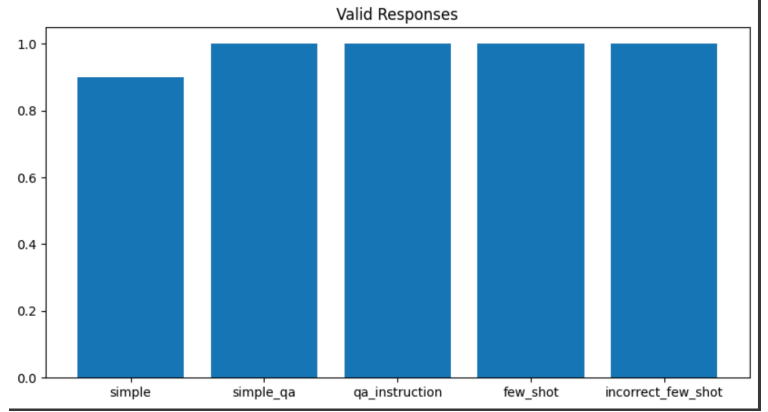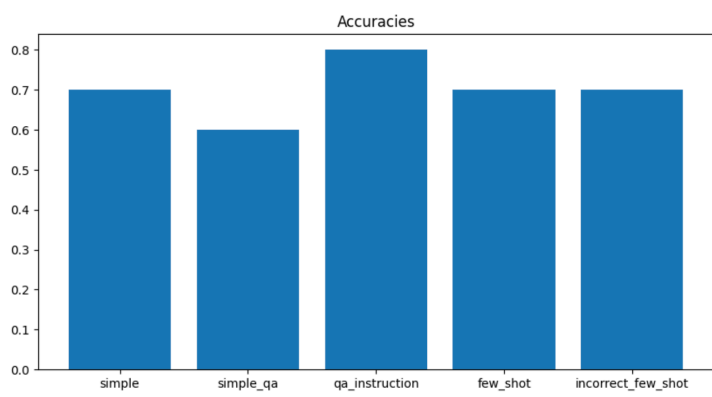
It's the crucial to provide protection to biodiversity and encourage environmental attention. I may assist you respond to some animal and plant associated inquiries, assist in conservation efforts, or present more data about their replica and critical requirements.

(iv)   Questions regarding personal chores for which the model has not seen the data.

(v)   Asking model to generate LinkedIn Biography while sequentially feeding it more and more information about one while restricting the number of words as output
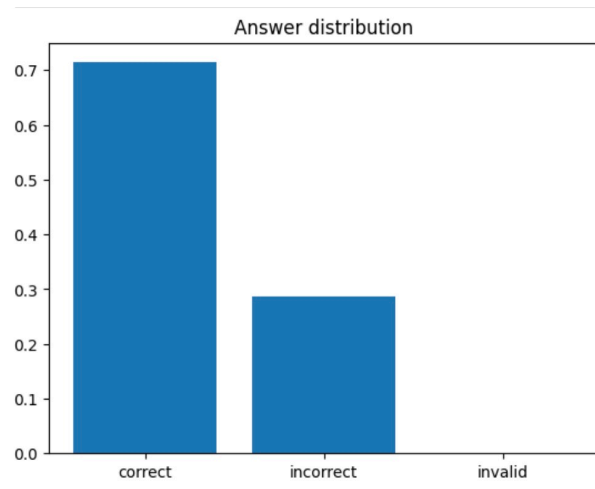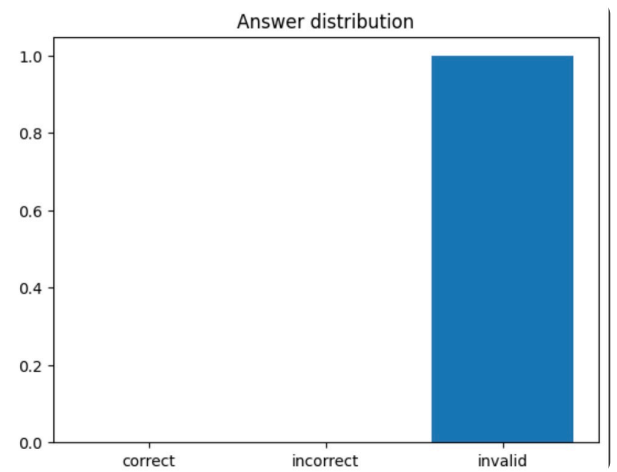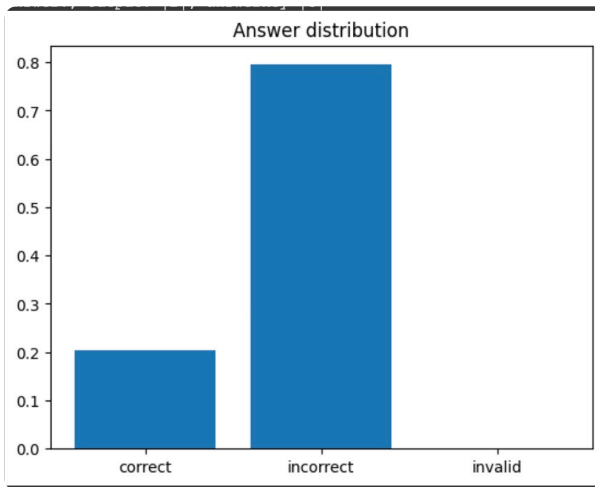
(b) Using LMs for classification :-

Prompt types :
(a) Simple prompt
(b) Simple QA prompt
(c) QA Instruction prompt
(d) Few shot Prompt
(e) Incorrect Few shot Prompt

(i)   Model is either predicting the wrong answer choices or invalid solutions.

(ii)   Similar accuracy with incorrect labels in the prompts as with correct and clean prompts.

(iii)   Model is more confident when it is correct.

Accuracies

Valid Responses

Confidence (Probability of most likely token)

(iv) GPT-2 Model has fewer parameters

(v)

Answer distribution

Answer distribution

Answer distribution

# Q3    Self-Prompting    Language Models

**(a)**

```
["Tokens" 1-5: soft prompt] [Tokens 6-50: question]
[Tokens 51-70: chain of thought reasoning]
[Token 71: answer] [Token 72: newline]
[Tokens 73-100: padding].
```

included tokens ⇒ 50 - 71  ⇒ we need the model to becom giving a chain of thought (50 - 69) as well as the answer (70) as well as ending with a new line (71)

**(b)** number of params trained in the model?

max sequence length ⇒ $S$
token embedding ⇒ $E$ →
Vocab size ⇒ $V$ →
Hidden State ⇒ $H$
Number of layers ⇒ $L$
Attention Query / Key dim ⇒ $D$

since we have 5 soft-prompt tokens that we are learning

$\underline{5E \; Params}$

**(c)** **(i)** TRUE

Once representations are computed during an initial forward pass, they can be stored & reused since once the representations are computed, they become context independent

**(ii)** TRUE

Best possible hard prompt $\subseteq$ soft prompts

**(iii)** FALSE

overfitting could occur or less of generalisation or catastrophic forgetting

**(iv)** FALSE

since only prompt embeddings are learnt. Rest of the model remains frozen.

(d) Describe how you would adapt a meta-learning approach like MAML for this situation?

Via MAML, the goal is to learn an initial condition of the prompt.

For each family of tasks, break the dataset into train, validation, and test. Initialise soft-prompt to current initialisation (will be random initially). Run SGD with $k$-steps to get $\phi_k$. Backprop to $\theta_0$. Update initialisation & repeat for another task

## Q4 : <u>Meta-learning for Learning 1D functions</u> :

<u>toy example</u> :

$\mathbb{R} \to \mathbb{R}$ $\to$ considering functions that map a scalar onto another.

we pick, $d$ functions $\phi_0, \phi_1, \dots \phi_{d-1}$

$$f_\alpha = \sum_{k=0}^{d-1} \alpha_k \phi_k(x)$$

Learner :

$$f_{\hat{\beta}, c} = \sum_{k=0}^{d-1} \hat{B}_k C_k \phi_k(x)$$

where $\hat{B} = \arg\min_{B} \|B\|_2^2$

s.b. $y_i = \sum_{k=0}^{d-1} B_k C_k \phi_k(x_i) \quad \forall i = 1, 2, \dots K$

(a) $(x, y)$, $\phi_t(x) = \phi_1(x)$

$$\phi_q(x) = \phi_2(x) = \phi_1(x)$$

$$y = \phi_t(x) = \phi_1(x) = \phi_2(x) \Rightarrow \phi(x)$$

for min-norm solution,

$$\hat{B} = A^T (AA^T)^{-1} y \qquad \text{where} \qquad A = \begin{bmatrix} c_0 \phi_0(x) & c_1 \phi_1(x) \end{bmatrix}$$

So,

$$\hat{B} = \phi(x) \begin{bmatrix} c_0 \\ c_1 \end{bmatrix} \cdot \left( (c_0^2 + c_1^2) \phi(x)^2 \right)^{-1} \cdot \phi(x)$$

$$\hat{B} = \begin{bmatrix} c_0 \\ c_1 \end{bmatrix} \cdot \frac{1}{c_0^2 + c_1^2}$$

(b) Calculate the gradient of the expected test error w.r.t. the feature weights $c_0$ and $c_1$ respectively:

$$\frac{\partial}{\partial c} \left( E_{x_{test}, y_{test}} \left\{ \frac{1}{2} \left\| y - \hat{B}_0 c_0 \phi_t(x) - \hat{B}_1 c_1 \phi_a(x) \right\|_2^2 \right) \right.$$

$$\to \phi_t(x) \qquad \phi_a(x)$$

$$E = E_{(x_{test}, y_{test})} \left[ \frac{1}{2} \left\| \phi_t(x) - \frac{c_0^2}{c_0^2 + c_1^2} \phi_t(x) - \frac{c_1^2}{c_0^2 + c_1^2} \phi_a(x) \right\|_2^2 \right.$$

Due to orthognality, $E \| X \|_2^2 = E[X_1^2 + X_2^2]$

$$= \frac{1}{2} \left( 1 - \frac{c_0^2}{c_0^2 + c_1^2} \right)^2 + \frac{1}{2} \left( \frac{c_1^2}{c_0^2 + c_1^2} \right)$$

(d) Using meta-learned feature weights the test performance is better than all the 1s feature weights but worse than the oracle.

For increasing $n$, the test error goes down.

For $n = 32$, the downward spike occurs since we are using $n\_train\_inner = 32$ for inward loop.

e) Weight an feature of interest grows with the process

f) for both num_gd-steps, performance of the model does improve.

_____