

Speculative Decoding

正田 備也

masada@rikkyo.ac.jp

Decoding の効率化

- ▶ autoregressive LLM は、次の 1 トークンの生成を繰り返すことで、長いテキストを生成する。
- ▶ しかし、そのための計算コストは大きい。
- ▶ そこで、小規模な LM を補助的に使って、トークンをより高速に生成する。具体的には・・・
- ▶ speculative decoding を使うと、元の LLM からトークン列を生成するときに使われるのと全く同じ確率分布に基づきつつも、より効率的にトークン列を生成できるようになる。

Speculative decoding [arXiv:2211.17192][arXiv:2302.01318]

- ▶ 二つのカテゴリカル分布を使ってトークンを生成する。
 - ▶ $\text{Cat}(\theta)$: 小規模 LM のある position でのカテゴリカル分布
 - ▶ $\text{Cat}(\phi)$: LLM の同じ position でのカテゴリカル分布
- 1. $\text{Cat}(\theta)$ からサンプリングし、得られたアイテムを x とする。
- 2. $\theta_x \leq \phi_x$ ならば x をそのまま採用する。
- 3. $\theta_x > \phi_x$ ならば、表の出る確率が $\frac{\phi_x}{\theta_x}$ のコインを投げる。
- 4. 表が出たら x を採用する。
- 5. 裏が出たら $\text{Cat}((\phi - \theta)_+)$ からサンプリングし直す。

サンプリングし直すときに使うカテゴリカル分布

- ▶ $\text{Cat}((\phi - \theta)_+)$ は、アイテム x の出現確率が、以下のように定義されるカテゴリカル分布である。

$$(\phi_x - \theta_x)_+ \equiv \frac{\max(0, \phi_x - \theta_x)}{\sum_{x'} \max(0, \phi_{x'} - \theta_{x'})} \quad (1)$$

- ▶ つまり、 $\phi_x > \theta_x$ を満たすアイテム x だけが非ゼロの確率を持つように規格化されたカテゴリカル分布である。

Speculative decoding におけるアイテムの出現確率

- ▶ 上の rejection sampling で結果的に選ばれるアイテムを確率変数 X で表すと、 $X = x$ となる確率は以下のように書ける。

$$\begin{aligned} P(X = x) &= P_{\theta}(X = x)P(x \text{ が採用される} | X = x) \\ &\quad + \left(P_{\phi}(X = x | \text{Cat}(\theta) \text{ からのサンプルが採用されない}) \right. \\ &\quad \left. \times P(\text{Cat}(\theta) \text{ からのサンプルが採用されない}) \right) \end{aligned} \quad (2)$$

- ▶ この式を、以下、簡単な形に書き換える。

まず、

$$P(x \text{ が採用される} | X = x) = \begin{cases} 1 & \text{when } \theta_x \leq \phi_x \\ \frac{\phi_x}{\theta_x} & \text{when } \theta_x > \phi_x \end{cases} \quad (3)$$

この確率は、以下のように書き換えることができる。

$$P(x \text{ が採用される} | X = x) = \min \left(1, \frac{\phi_x}{\theta_x} \right) \quad (4)$$

$P_{\theta}(X = x) \equiv \theta_x$ であるから、式 (2) の一つ目の項は、以下のように書き換えられる。

$$P_{\theta}(X = x)P(x \text{ が採用される} | X = x) = \theta_x \times \min \left(1, \frac{\phi_x}{\theta_x} \right) = \min(\theta_x, \phi_x) \quad (5)$$

次に、式 (5) より

$$P_{\theta}(X' = x')P(x' \text{ が採用されない} | X' = x') = \theta_{x'} \times \left(1 - \min\left(1, \frac{\phi_{x'}}{\theta_{x'}}\right)\right) \quad (6)$$

ところで、 $\theta_x \leq \phi_x$ のとき $1 - \min\left(1, \frac{\phi_{x'}}{\theta_{x'}}\right) = 1 - 1 = 0$ であり、また、 $\theta_x > \phi_x$ のとき $1 - \min\left(1, \frac{\phi_{x'}}{\theta_{x'}}\right) = 1 - \frac{\phi_{x'}}{\theta_{x'}}$ であるので、両方をまとめると、

$$1 - \min\left(1, \frac{\phi_{x'}}{\theta_{x'}}\right) = \max\left(0, 1 - \frac{\phi_{x'}}{\theta_{x'}}\right) \quad (7)$$

と書ける。したがって、

$$\begin{aligned} P_{\theta}(X' = x')P(x' \text{ が採用されない} | X' = x') &= \theta_{x'} \times \left(1 - \min\left(1, \frac{\phi_{x'}}{\theta_{x'}}\right)\right) \\ &= \theta_{x'} \times \max\left(0, 1 - \frac{\phi_{x'}}{\theta_{x'}}\right) \\ &= \max(0, \theta_{x'} - \phi_{x'}) \end{aligned} \quad (8)$$

式 (8) を使うと、 $\text{Cat}(\theta)$ からのサンプルが採用されない確率は、

$$\begin{aligned} P(\text{Cat}(\theta) \text{ からのサンプルが採用されない}) &= \sum_{x'} P(X' = x', x' \text{ が採用されない}) \\ &= \sum_{x'} P_{\theta}(X' = x') P(x' \text{ が採用されない} | X' = x') \\ &= \sum_{x'} \max(0, \theta_{x'} - \phi_{x'}) \end{aligned} \tag{9}$$

と書けることが分かる。

ここで、 $\text{Cat}(\theta)$ からのサンプルが採用されないときには、アイテム x の出現確率が

$$(\phi_x - \theta_x)_+ \equiv \frac{\max(0, \phi_x - \theta_x)}{\sum_{x'} \max(0, \phi_{x'} - \theta_{x'})} \tag{10}$$

と定義されるカテゴリカル分布からサンプリングするのだったことを、思い出す。

式 (10) の分母は、式 (9) に一致している。よって、式 (2) の二つ目の項は以下のようになる。

$$\begin{aligned} & P_{\phi}(X = x | \text{Cat}(\theta) \text{ からのサンプルが採用されない}) \\ & \times P(\text{Cat}(\theta) \text{ からのサンプルが採用されない}) \\ & = (\phi_x - \theta_x)_+ \times \sum_{x'} \max(0, \theta_{x'} - \phi_{x'}) \\ & = \frac{\max(0, \phi_x - \theta_x)}{\sum_{x'} \max(0, \phi_{x'} - \theta_{x'})} \times \sum_{x'} \max(0, \theta_{x'} - \phi_{x'}) \\ & = \max(0, \phi_x - \theta_x) \end{aligned} \tag{11}$$

以上より、式 (2) は、以下のように書き換えられる。

$$P(X = x) = \min(\theta_x, \phi_x) + \max(0, \phi_x - \theta_x) \tag{12}$$

θ_x と ϕ_x の大小関係で場合分けして計算すると、

$$P(X = x) = \phi_x \tag{13}$$