

大模型面经之解码策略topk、topp（附代码）

瓦力算法学研所 2024年05月13日 18:41 北京

介绍

大模型文本生成中的解码策略是用来控制生成文本的多样性和连贯性的方法。当模型为下一个词生成概率分布时，这些策略决定了从概率分布中选择词的方式。

常用的解码方式有：

- 贪心解码（Greedy Decoding）：直接选择概率最高的单词。这种方法简单高效，但是可能会导致生成的文本过于单调和重复。
- 随机采样（Random Sampling）：按照概率分布随机选择一个单词。这种方法可以增加生成的多样性，但是可能会导致生成的文本不连贯和无意义。
- Beam Search：维护一个大小为 k 的候选序列集合，每一步从每个候选序列的概率分布中选择概率最高的 k 个单词，然后保留总概率最高的 k 个候选序列。这种方法可以平衡生成的质量和多样性，但是可能会导致生成的文本过于保守和不自然

Top-k、Top-p和Temperature是三种LLM常用的解码策略

Top-k

Top-k采样是一种简单的截断策略，它从一个限制为前k个最可能的词的集合中随机挑选下一个词。"k"是一个可以调整的参数，较小的k值倾向于生成更可预测的文本，而较大的k值则鼓励多样性。当k很大时，几乎等同于无约束采样；当k=1时，该策略退化为贪心解码，即永远选择最可能的下一个词。

```
1 import torch
2 from labml_nn.sampling import Sampler
3
4 # Top-k Sampler
5 class TopKSampler(Sampler):
6     # k is the number of tokens to pick
7     # sampler is the sampler to use for the top-k tokens
8     # sampler can be any sampler that takes a logits tensor as input and returns a tensor of token indices
9     def __init__(self, k: int, sampler: Sampler):
10         self.k = k
11         self.sampler = sampler
12
13     # Sample from logits
14     def __call__(self, logits: torch.Tensor):
15         # New logits filled with -∞; i.e. zero probability
```

```

16     zeros = logits.new_ones(logits.shape) * float('-inf')
17     # Pick the largest k logits and their indices
18     values, indices = torch.topk(logits, self.k, dim=-1)
19     # Set the values of the top-k selected indices to actual logits.
20     # Logits of other tokens remain -∞
21     zeros.scatter_(-1, indices, values)
22     # Sample from the top-k logits with the specified sampler.
23     return self.sampler(zeros)

```

Top-p (Nucleus) Sampling

Top-p采样，又称为nucleus采样，是另一种更加动态的截断策略。它先对所有可能的下一个词的概率进行排序，**然后累加这些概率，直到达到预先设定的阈值p**。这样选择的词集合将包含累计概率至少为p的集合，这通常意味着包含多数概率大词但只包含部分词汇。这样，生成的文本既能够保持多样性，又能保持一定的连贯性。

```

1  import torch
2  from torch import nn
3
4  from labml_nn.sampling import Sampler
5
6
7  class NucleusSampler(Sampler):
8      """
9      ## Nucleus Sampler
10     """
11     def __init__(self, p: float, sampler: Sampler):
12         """
13         :param p: is the sum of probabilities of tokens to pick $p$
14         :param sampler: is the sampler to use for the selected tokens
15         """
16         self.p = p
17         self.sampler = sampler
18         # Softmax to compute $P(x_i | x_{1:i-1})$ from the logits
19         self.softmax = nn.Softmax(dim=-1)
20
21     def __call__(self, logits: torch.Tensor):
22         """
23         Sample from logits with Nucleus Sampling
24         """
25
26         # Get probabilities $P(x_i | x_{1:i-1})$

```

```

27     probs = self.softmax(logits)
28
29     # Sort probabilities in descending order
30     sorted_probs, indices = torch.sort(probs, dim=-1, descending=True)
31
32     # Get the cumulative sum of probabilities in the sorted order
33     cum_sum_probs = torch.cumsum(sorted_probs, dim=-1)
34
35     # Find the cumulative sums less than $p$.
36     nucleus = cum_sum_probs < self.p
37
38     # Prepend ones so that we add one token after the minimum number
39     # of tokens with cumulative probability less than $p$.
40     nucleus = torch.cat([nucleus.new_ones(nucleus.shape[:-1] + (1,)), nucleus], dim=-1)
41
42     # Get log probabilities and mask out the non-nucleus
43     sorted_log_probs = torch.log(sorted_probs)
44     sorted_log_probs[~nucleus] = float('-inf')
45
46     # Sample from the sampler
47     sampled_sorted_indexes = self.sampler(sorted_log_probs)
48
49     # Get the actual indexes
50     res = indices.gather(-1, sampled_sorted_indexes.unsqueeze(-1))
51
52     #
53     return res.squeeze(-1)

```

Temperature

Temperature用来平滑或加剧模型的概率分布。温度参数T影响每个词被选中的概率。通过改变概率分布的锐度，可以控制生成的文本的随机性：

$$\rho_i = \frac{1}{Q} e^{-\epsilon_i/kT} = \frac{e^{-\epsilon_i/kT}}{\sum_{j=1}^M e^{-\epsilon_j/kT}}$$

当 $T > 1$ 时，分布变得更平缓，小概率事件被选中的机会增大，这增加了生成文本的随机性和多样性。

当 $T = 1$ 时，保留模型原始的概率分布。

当 $T < 1$ 时，分布变得更尖锐，大概率事件更有可能被选中，这使得文本生成更加确定性，并且倾向于重复性。

```
1 import torch
2 from torch.distributions import Categorical
3
4 from labml_nn.sampling import Sampler
5
6
7 class TemperatureSampler(Sampler):
8     """
9     ## Sampler with Temperature
10    """
11    def __init__(self, temperature: float = 1.0):
12        """
13        :param temperature: is the temperature to sample with
14        """
15        self.temperature = temperature
16
17    def __call__(self, logits: torch.Tensor):
18        """
19        Sample from logits
20        """
21
22        # Create a categorical distribution with temperature adjusted logits
23        dist = Categorical(logits=logits / self.temperature)
24
25        # Sample
26        return dist.sample()
```

配合使用这些策略可以极大影响生成文本的品质。实践中，人们经常根据具体应用调整这些参数，以便获得既丰富多样又逻辑连贯的文本。实际上，这些策略也经常被一起使用。例如，可以同时使用top-k和top-p来过滤候选词汇（topk > topp, 也即先选取topk个最大概率的候选项，再从这些候选项中概率和满足p的前n个选项），再通过调整temperature来进行最终的词选取。

文章来源：

https://www.zhihu.com/tardis/bd/art/647813179?source_id=1001

本系列将会持续更新，想要获取面经资料的同学欢迎关注公众号，进群一起交流~