

原理概要

在 GPT-2 中，**Dropout** 是一种重要的正则化技术，主要用于**防止模型在训练过程中过拟合**。它通过在训练时随机“关闭”（即将其输出设置为0）部分神经元来工作。

Dropout 层的作用

Dropout 层的核心作用是**提高模型的泛化能力**。在训练期间，它会随机地将网络中某些神经元的输出设置为零。这会带来几个关键好处：

1. **防止神经元协同适应**：如果没有 dropout，神经元可能会依赖于其他特定神经元的存在，形成一种“协同适应”的关系。这导致模型对训练数据中的特定模式过于敏感，但对新数据却表现不佳。通过随机关闭神经元，dropout 强制网络中的每个神经元都必须独立学习有用的特征，使其不能过度依赖其他神经元。
2. **模拟集成学习**：每次进行前向传播时，dropout 都会生成一个略有不同的网络“子集”。这相当于在训练过程中训练了大量不同的、更小的模型。在推理阶段，所有神经元都被激活，这可以被看作是这些“子模型”的预测结果的平均，从而获得更稳健、更泛化的结果。

重要的是，**dropout 只在训练阶段使用**。在推理（或评估）阶段，所有神经元都保持激活，因为我们希望使用整个模型的全部能力来进行预测。

dropout_rate 和 scale 的作用

在 `torch.nn.Dropout` 或类似的实现中，`dropout_rate` 和 `scale` 是两个关键参数。

dropout_rate

`dropout_rate`（通常也写作 `p`）是一个浮点数，其值介于 0 和 1 之间。它代表了在训练过程中，每个神经元被设置为零的概率。

- 例如，如果 `dropout_rate = 0.1`，这意味着在每个训练步骤中，输入张量中的约 10% 的元素会被随机地设置为0。

选择合适的 `dropout_rate` 至关重要。过高的 `dropout_rate` 可能会导致模型欠拟合，因为它丢弃了太多有用的信息；而过低的 `dropout_rate` 则可能无法有效防止过拟合。在 GPT-2 中，常见的 `dropout_rate` 设置是 0.1。

scale（缩放）

`scale` 是与 `dropout` 机制紧密相关的操作。当一部分神经元被随机丢弃后，剩下的激活神经元需要进行缩放，以保持网络总体的激活值不变。

其缩放因子为 $1 / (1 - \text{dropout_rate})$ 。

这个缩放操作确保了在训练和推理阶段，神经元的期望输出保持一致。

为什么需要缩放？

假设一个神经元在训练时有 50% 的概率被丢弃（`dropout_rate = 0.5`）。这意味着它的期望输出会减半。如果在推理时所有神经元都激活，那么这个神经元的输出将是训练时的两倍，这会导致网络的激活值发生变化，从而影响模型的性能。通过将未丢弃的神经元输出乘以 $1 / (1 - 0.5) = 2$ ，我们可以确保：

- **训练时：** $0.5 * 0$ （丢弃的） + $0.5 * (\text{输出} * 2)$ （保留的） = 输出
- **推理时：** $1 * \text{输出}$

可以看到，无论是在训练还是推理时，这个神经元的期望输出都是一致的，从而消除了训练和推理之间的不一致性。在 PyTorch 等框架中，这种缩放是**自动完成的**。

前向传播

首先我们要生成一个与输入张量形状相同的掩码张量，并且用[0,1]之间的均匀分布来初始化这个张量，我们还需要缓存这个张量，以便于反向传播时使用。

将经过初始化的掩码张量作用于输入张量，这是一个逐元素比较过程。如果当前掩码元素的值大于丢弃率 (`dropout_rate`)，则将掩码张量当前元素设置为 1，否则设置为 0，其次再乘以输入张量的对应元素，同时乘以缩放因子 (`scale`)。经过这个操作之后，原始的输入张量的某些地方变成了 0，其余的地方被乘以缩放因子。

反向传播

假设输入梯度已知，其形状与前向传播的输入一致。

我们将前向传播的掩码张量作用于输入梯度，这也是一个逐元素操作过程，如果掩码值为1，则输入梯度乘以缩放因子，否则梯度设置为零。