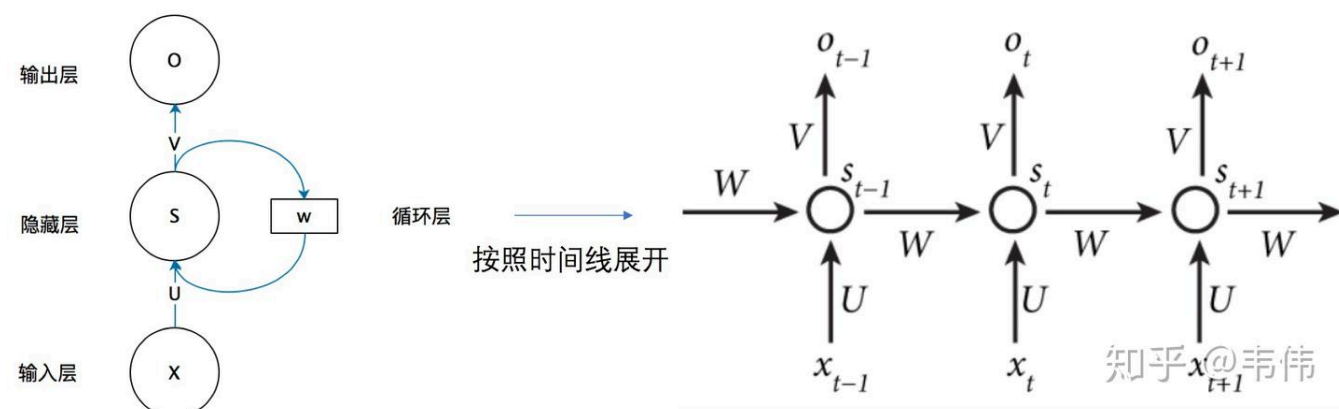
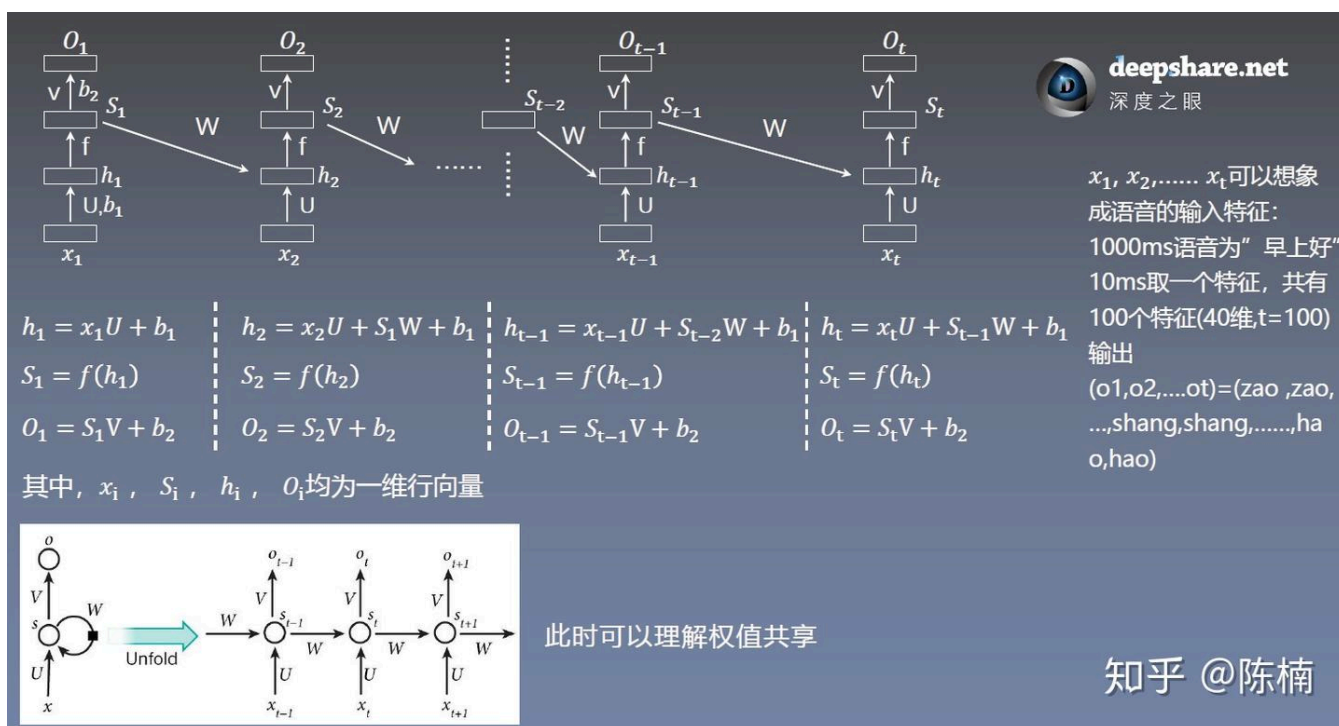


# RNN,LSTM,GRU计算方式

## RNN计算方式

```
torch.nn.RNN(input_size, hidden_size, num_layers=1, nonlinearity='tanh', bias=True,
batch_first=False, dropout=0.0, bidirectional=False, device=None, dtype=None)
```



$$h_t = \tanh(x_t W_{ih}^T + b_{ih} + h_{t-1} W_{hh}^T + b_{hh})$$

```
# Efficient implementation equivalent to the following with bidirectional=False
def forward(x, h_0=None):
    if batch_first:
        x = x.transpose(0, 1)
    seq_len, batch_size, _ = x.size()
    if h_0 is None:
        h_0 = torch.zeros(num_layers, batch_size, hidden_size)
    h_t_minus_1 = h_0
    h_t = h_0
    output = []
    for t in range(seq_len):
```

```

    for layer in range(num_layers):
        h_t[layer] = torch.tanh(
            x[t] @ weight_ih[layer].T
            + bias_ih[layer]
            + h_t_minus_1[layer] @ weight_hh[layer].T
            + bias_hh[layer]
        )
    output.append(h_t[-1])
    h_t_minus_1 = h_t
output = torch.stack(output)
if batch_first:
    output = output.transpose(0, 1)
return output, h_t

```

- 参数：
  - input\_size-输入X中预期特征的数量
  - hidden size-隐藏状态h中的特征数量
  - num\_layers-循环层数。例如，设置num\_layers:=2表示将两个RNN堆叠在一起形成一个堆叠RNN，第二个RNN接收第一个RNN的输出并计算最终结果。默认为：1
  - nonlinearity-使用的非线性函数。可以是'tanh'或'relu'。默认为："tanh"
  - bias-如果为False,则该层不使用偏置权重bh和bhh。默认为：True
  - batch\_first-如果为True,则输入和输出张量提供为(batch,seq,feature),而不是(seg,batch,feature)。请注意，这不适用于隐藏状态或单元状态。有关详细信息，请参阅下面的输入/输出部分。默认为：False
  - dropout-如果非零，则在除最后一层外的每个RNN层的输出上引入Dropout层，dropout概率等于dropout。默认为：0
  - bidirectional-如果为True,则成为双向RNN。默认为：False

输入：input,hx

·input:形状为(L,Hin)的张量（用于无批次输入），形状为(L,N,Hn)当batch\_first:=False时，或形状为(N,L,Hin)当batch\_first=True时，包含输入序列的特征。输入也可以是填充的可变长度序列。有关详细信息，请参阅torch.nn.utils.rnn.pack\_padded\_sequence()或torch.nn.utils.rnn.pack\_sequence()

·hx:形状为(Dnum\_layers,Hout)的张量（用于无批次输入），或形状为(Dnum\_layers,.N,Hout)的张量（用于输入序列批次）包含初始隐藏状态。如果未提供，则默认为零。

其中

$$\begin{aligned}
 N &= \text{batch size} \\
 L &= \text{sequence length} \\
 D &= 2 \text{ if } \text{bidirectional}=\text{True} \text{ otherwise } 1 \\
 H_{in} &= \text{input size} \\
 H_{out} &= \text{hidden size}
 \end{aligned}$$

输出：output,hn

output:形状为(L,DHout)的张量（用于无批次输入），或形状为(L,N,DHout)当batch\_first:=False时，或形状为(W,L,D\*Hout)当batch\_first=True时，包含从RNN最后一层在每个t的输出特征ht)。如果输入是torch.nn.utils.rnn.PackedSequence,则输出也将是打包序列。

hn:形状为(Dnum\_layers,Hout)的张量（用于无批次输入），或形状为Dnum\_layers,N,Hout)的张量（用于批次中的每个元素）包含最终隐藏状态。

变量：

weight ih[k]-第k层的可学习输入-隐藏权重，对于k=0,形状为hidden size,input size)。否则，形状(hidden size,num\_directions hidden\_size)

weight\_hh[k]-第k层的可学习隐藏-隐藏权重，形状为hidden\_size,.hidden\_size)

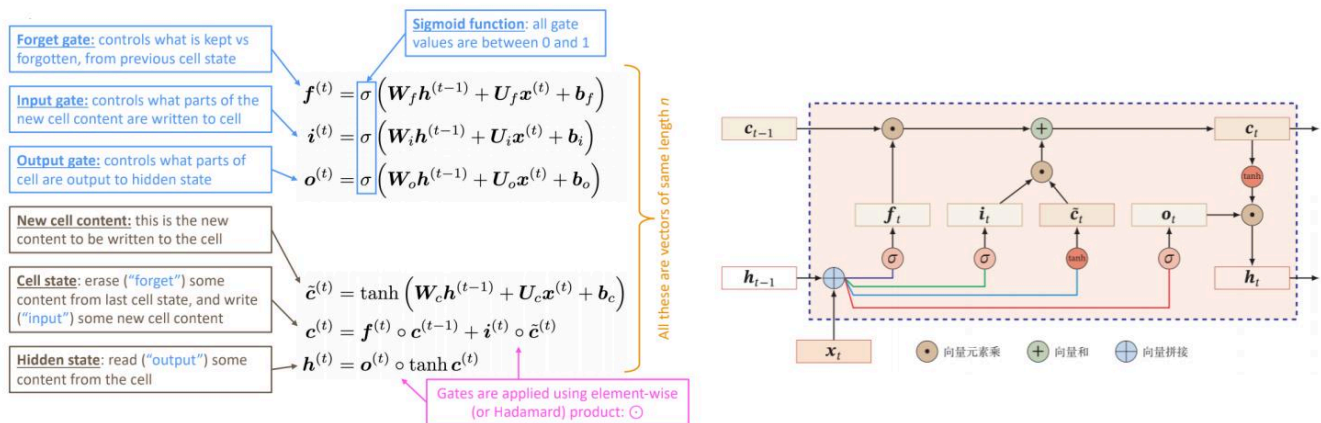
bias ih I[k]-第k层的可学习输入-隐藏偏置，形状为hidden\_size)

bias\_hh[k]-第k层的可学习隐藏-隐藏偏置，形状为hidden\_size)

```
rnn = nn.RNN(10, 20, 2)
input = torch.randn(5, 3, 10)
h0 = torch.randn(2, 3, 20)
output, hn = rnn(input, h0)
```

## LSTM计算方式

```
torch.nn.LSTM(input_size, hidden_size, num_layers=1, bias=True, batch_first=False, dropout=0.0,
bidirectional=False, proj_size=0, device=None, dtype=None)
```



$$\begin{aligned}i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{t-1} + b_{hi}) \\f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}h_{t-1} + b_{hf}) \\g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{t-1} + b_{hg}) \\o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}h_{t-1} + b_{ho}) \\c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\h_t &= o_t \odot \tanh(c_t)\end{aligned}$$

在时间步 t,  $h_t$  表示隐藏状态,  $C_t$  表示细胞状态,  $x_t$  表示时间步 t 的输入,  $h_{t-1}$  表示上一时间步 (t-1) 的隐藏状态, 或在初始时刻 t=0 时的初始隐藏状态。it、ft、gt、ot 分别表示输入门、遗忘门、细胞门和输出门。 $\sigma$  是 sigmoid 激活函数,  $\odot$  表示哈达玛积 (逐元素相乘)。

在多层 LSTM 中, 第 l 层的输入是前一层隐藏状态乘以一个 dropout 掩码, 其中每个元素是一个伯努利随机变量, 以概率 p 为 0 (即 dropout 概率为 p)。

如果指定了 proj\_size > 0, 则将使用带有投影的 LSTM (LSTM with projections), 这会对 LSTM 单元做如下修改:

第一,  $h_t$  的维度将从 hidden\_size 改变为 proj\_size (相应地, 权重矩阵  $W_{hi}$  的维度也会改变)。

第二, 每一层输出的隐藏状态将乘以一个可学习的投影矩阵:  $h_t = W_{hr} \cdot h_t$ , 其中  $W_{hr}$  是投影矩阵。

需要注意的是, 由于这一修改, LSTM 网络的输出形状也会随之改变。关于所有变量的精确维度, 请参见下文的“输入/输出”部分。

更多细节可参考论文: <https://arxiv.org/abs/1402.1128>。

参数:

`input_size`-输入 $x$ 中预期特征的数量

`hidden size`-隐藏状态 $h$ 中的特征数量·`num_layers`-循环层数。例如, 设置`num_layers=2`意味着将两个LSTM堆叠在一起形成一个“堆叠 LSTM”,第二个LSTM接收第一个LSTM的输出来计算最终结果。默认为1。

`bias`-如果为False,则该层不使用偏置权重 $b_h$ 和 $b_{hh}$ 。默认为True

`batch_first`-如果为True,则输入和输出张量将提供为(batch,seq,feature),而不是(seq,batch,feature)o请注意, 这不适用于隐藏状态或单元状态。有关详细信息, 请参阅下面的“输入/输出”部分。默认为False。

`dropout`-如果非零, 则在除最后一层外的每个LSTM层的输出上引入Dropout层, dropout概率等于 dropout。默认为0。

`bidirectional`-如果为True,则变为双向LSTM。默认为False。

`proj_size`-如果 $>0$ , 则将使用具有相应大小的投影的LSTM。默认为0。

输入: `input,(h0,c0)`

`input`:形状为 $(L, H_n)$ 的张量, 用于未批处理的输入; 当`batch_first=False`时, 形状为 $(L, N, H_m)$ , 或者当`batch_first=True`时, 形状为 $(N, L, H_m)$ , 其中包含输入序列的特征。输入也可以是打包的可变长度序列。有关详细信息, 请参阅`torch.nn.utils.rnn.pack_padded_sequence()`或`torch.nn.utils.rnn.pack_sequence()`

`h0`:形状为 $(D_{num\_layers}, H_{ot})$ 的张量, 用于未批处理的输入; 当`batch first=False`时, 形状为 $(D_{num\_layers}, N, H_{out})$ 的张量, 用于批处理的输入, 包含输入序列的初始隐藏状态。如果未提供 `h0,c0`, 则默认为零。

`c_0`:形状为 $(D_{num\_layers}, H_{ce})$ 的张量, 用于未批处理的输入; 当`batch_first=False`时, 形状为 $(D_{num\_layers}, N, H_{cen})$ 的张量, 用于批处理的输入, 包含输入序列的初始单元状态。如果未提供 `h0,c0`, 则默认为零

其中

$$N = \text{batch size}$$
$$L = \text{sequence length}$$
$$D = 2 \text{ if } \text{bidirectional}=\text{True} \text{ otherwise } 1$$
$$H_{in} = \text{input}_{\text{size}}$$
$$H_{cell} = \text{hidden}_{\text{size}}$$
$$H_{out} = \text{proj}_{\text{size}} \text{ if } \text{proj}_{\text{size}} > 0 \text{ otherwise } \text{hidden}_{\text{size}}$$

输出: `output,(hn,cn)`

`output`:形状为 $(L, D_{Hot})$ 的张量, 用于未批处理的输入; 当`batch first=False`时, 形状为 $(L, N, D_{Hot})$ , 或者当`batch_first=True`时, 形状为 $(N, L, D * H_{out})$ , 包含LSTM最后一个时间步的输出特征 $h_t$ 。如果输入是`torch.nn.utils.rnn.PackedSequence`, 则输出也将是打包序列。当 `bidirectional=True`时, `output`将包含每个时间步的前向和后向隐藏状态的连接。

`hn`:形状为 $(D_{num\_layers}, H_{ot})$ 的张量, 用于未批处理的输入; 当`batch_first=False`时, 形状为 $(D_{num\_layers}, N, H_{ot})$ 的张量, 用于批处理的输入, 包含序列中每个元素的最终隐藏状态。当`bidirectional=True`时, `hn`将分别包含最终前向和后向隐藏状态的连接。

cn:形状为( $Dnum\_layers, Hce$ )的张量, 用于未批处理的输入; 当 $batch\_first=False$ 时, 形状为( $Dnum\_layers, N, Hceu$ )的张量, 用于批处理的输入, 包含序列中每个元素的最终单元状态。当 $bidirectional=True$ 时, Gn将分别包含最终前向和后向单元状态的连接。

参数:

weight\_ih\_l[k]: 第 k 层的可学习输入到隐藏层权重 (对应输入门、遗忘门、细胞门、输出门, 即  $W_{ii}, W_{if}, W_{ig}, W_{io}$ ), 其形状为:

当  $k=0$  时: ( $4 \times hidden\_size, input\_size$ );

- 当  $k>0$  时: ( $4 \times hidden\_size, num\_directions \times hidden\_size$ );
- 若指定了  $proj\_size>0$  且  $k>0$ : 形状为 ( $4 \times hidden\_size, num\_directions \times proj\_size$ )。

weight\_hh\_l[k]: 第 k 层的可学习隐藏层到隐藏层权重 (即  $W_{hi}, W_{hf}, W_{hg}, W_{ho}$ ), 其形状为:

- ( $4 \times hidden\_size, hidden\_size$ );
- 若指定了  $proj\_size>0$ , 则形状为 ( $4 \times hidden\_size, proj\_size$ )。

bias\_ih\_l[k]: 第 k 层的可学习输入到隐藏层偏置 (即  $b_{ii}, b_{if}, b_{ig}, b_{io}$ ), 形状为 ( $4 \times hidden\_size$ )。

bias\_hh\_l[k]: 第 k 层的可学习隐藏层到隐藏层偏置 (即  $b_{hi}, b_{hf}, b_{hg}, b_{ho}$ ), 形状为 ( $4 \times hidden\_size$ )。

weight\_hr\_l[k]: 第 k 层的可学习投影权重, 形状为 ( $proj\_size, hidden\_size$ )。

仅在指定了  $proj\_size>0$  时存在。

weight\_ih\_l[k]\_reverse: 反向传播方向中, 第 k 层的输入到隐藏层权重, 结构与 `weight_ih_l[k]` 类似。

仅在 `bidirectional=True` 时存在。

weight\_hh\_l[k]\_reverse: 反向传播方向中, 第 k 层的隐藏层到隐藏层权重, 结构与 `weight_hh_l[k]` 类似。

仅在 `bidirectional=True` 时存在。

bias\_ih\_l[k]\_reverse: 反向传播方向中, 第 k 层的输入到隐藏层偏置, 结构与 `bias_ih_l[k]` 类似。

仅在 `bidirectional=True` 时存在。

bias\_hh\_l[k]\_reverse: 反向传播方向中, 第 k 层的隐藏层到隐藏层偏置, 结构与 `bias_hh_l[k]` 类似。

仅在 `bidirectional=True` 时存在。

weight\_hr\_l[k]\_reverse: 反向传播方向中, 第 k 层的可学习投影权重, 结构与 `weight_hr_l[k]` 类似。

仅在 `bidirectional=True` 且指定了  $proj\_size>0$  时存在。

```

rnn = nn.LSTM(10, 20, 2)
input = torch.randn(5, 3, 10)
h0 = torch.randn(2, 3, 20)
c0 = torch.randn(2, 3, 20)
output, (hn, cn) = rnn(input, (h0, c0))

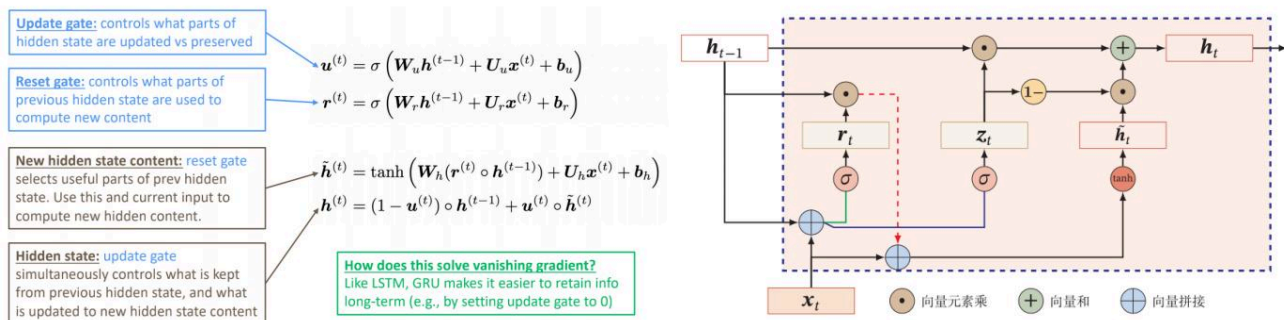
```

## GRU计算方式

```

torch.nn.GRU(input_size, hidden_size, num_layers=1, bias=True, batch_first=False, dropout=0.0,
bidirectional=False, device=None, dtype=None)

```



$$\begin{aligned}
 r_t &= \sigma(W_{ir}x_t + b_{ir} + W_{hr}h_{(t-1)} + b_{hr}) \\
 z_t &= \sigma(W_{iz}x_t + b_{iz} + W_{hz}h_{(t-1)} + b_{hz}) \\
 n_t &= \tanh(W_{in}x_t + b_{in} + r_t \odot (W_{hn}h_{(t-1)} + b_{hn})) \\
 h_t &= (1 - z_t) \odot n_t + z_t \odot h_{(t-1)}
 \end{aligned}$$

其中,  $h_t$ 是时间 $t$ 的隐藏状态,  $x_t$ 是时间 $t$ 的输入,  $h_{(t-1)}$ 是时间 $t-1$ 的层隐藏状态或时间0的初始隐藏状态,  $r_t$ 、 $z_t$ 、 $n_t$ 分别是重置门、更新门和新门。  $\sigma$ 是sigmoid函数,  $\odot$ 是阿达玛积(Hadamard product)。

在多层GRU中, 第 $l$ 层( $l \geq 2$ )的输入 $x_t$ 是前一层的隐藏状态 $h_{(t-1)}$ 乘以dropout, 其中每个dropout都是一个伯努利随机变量, 其取值为0的概率等于dropout。

参数:

input\_size-输入 $X$ 中预期特征的数量

hidden\_size-隐藏状态 $h$ 中的特征数量·num\_layers-循环层的数量。例如, 设置 num\_layers=2意味着将两个GRU堆叠在一起形成一个堆叠 GRU, 其中第二个GRU接收第一个GRU的输出并计算最终结果。默认值: 1

bias-如果为False, 则层不使用偏差权重 $b_h$ 和 $b_{hh}$ 。默认值: True

batch\_first-如果为True, 则输入和输出张量以batch, seq, feature形式提供, 而不是(seq, batch, feature)。请注意, 这不适用于隐藏状态或单元状态。有关详细信息, 请参阅下面的输入/输出部分。默认值: False

dropout-如果非零, 则在除最后一层之外的每个GRU层的输出上引入一个Dropout层, dropout概率等于 dropout。默认值: 0

bidirectional-如果为True, 则成为双向GRU。默认值: False

输入: input, h0



input:对于未批处理输入,形状为(L,Hm)的张量;当batch\_first=False时,形状为(L,W,Hn),当batch first=True时,形状为(N,L,Hn),包含输入序列的特征。输入也可以是填充后的可变长度序列。详情请参见 torch.nn.utils.rnn.pack\_padded\_sequence()或 torch.nn.utils.rnn.pack sequence()o。

h0:形状为(Dnum\_layers,H0t)或(Dnum layers,N,Hout)的张量,包含输入序列的初始隐藏状态。如果未提供,则默认为零。

其中:

$$\begin{aligned} N &= \text{batch size} \\ L &= \text{sequence length} \\ D &= 2 \text{ if bidirectional=True otherwise } 1 \\ H_{in} &= \text{input size} \\ H_{out} &= \text{hidden size} \end{aligned}$$

输出: output,h\_n

output:对于未批处理输入,形状为(L,DHout)的张量;当batch\_first=False时,形状为(L,N,DHout),当batch\_first=True时,形状为(N,L,DHout),包含GRU最后一层在每个t的输出特征(h\_t)。如果输入是 torch.nn.utils.rnn.PackedSequence,则输出也将是打包序列。

h\_n:形状为(Dnum\_layers,Hout)或(D\*num\_layers,N,Hout)的张量,包含输入序列的最终隐藏状态。

变量:

weight\_ih\_l[k-第kh层的可学习输入-隐藏权重(W\_irlW\_izlW\_in),当k=0时,形状为(3hidden\_size,input\_size)。否则,形状为(3hidden\_size,num\_directionshidden\_size)

weight\_hh\_l[-第kh层的可学习隐藏-隐藏权重(W\_hrW\_hzW\_hn),形状为(3 hidden\_size, hidden\_size)

bias\_ih\_l[-第k执层的可学习输入-隐藏偏差(b irb izb in),形状为(3 hidden\_size)

bias\_hh\_l[k-第kh层的可学习隐藏-隐藏偏差(b hrlb\_hzlb hn),形状为(3hidden\_size)

```
rnn = nn.GRU(10, 20, 2)
input = torch.randn(5, 3, 10)
h0 = torch.randn(2, 3, 20)
output, hn = rnn(input, h0)
```