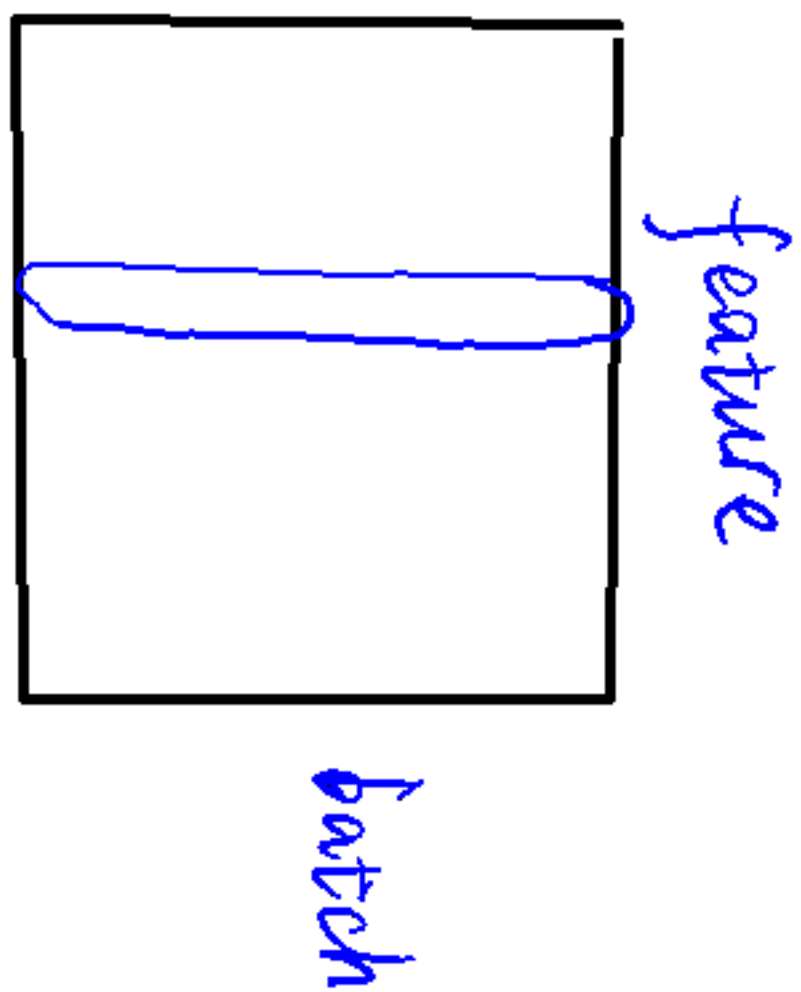


## Batch-Norm

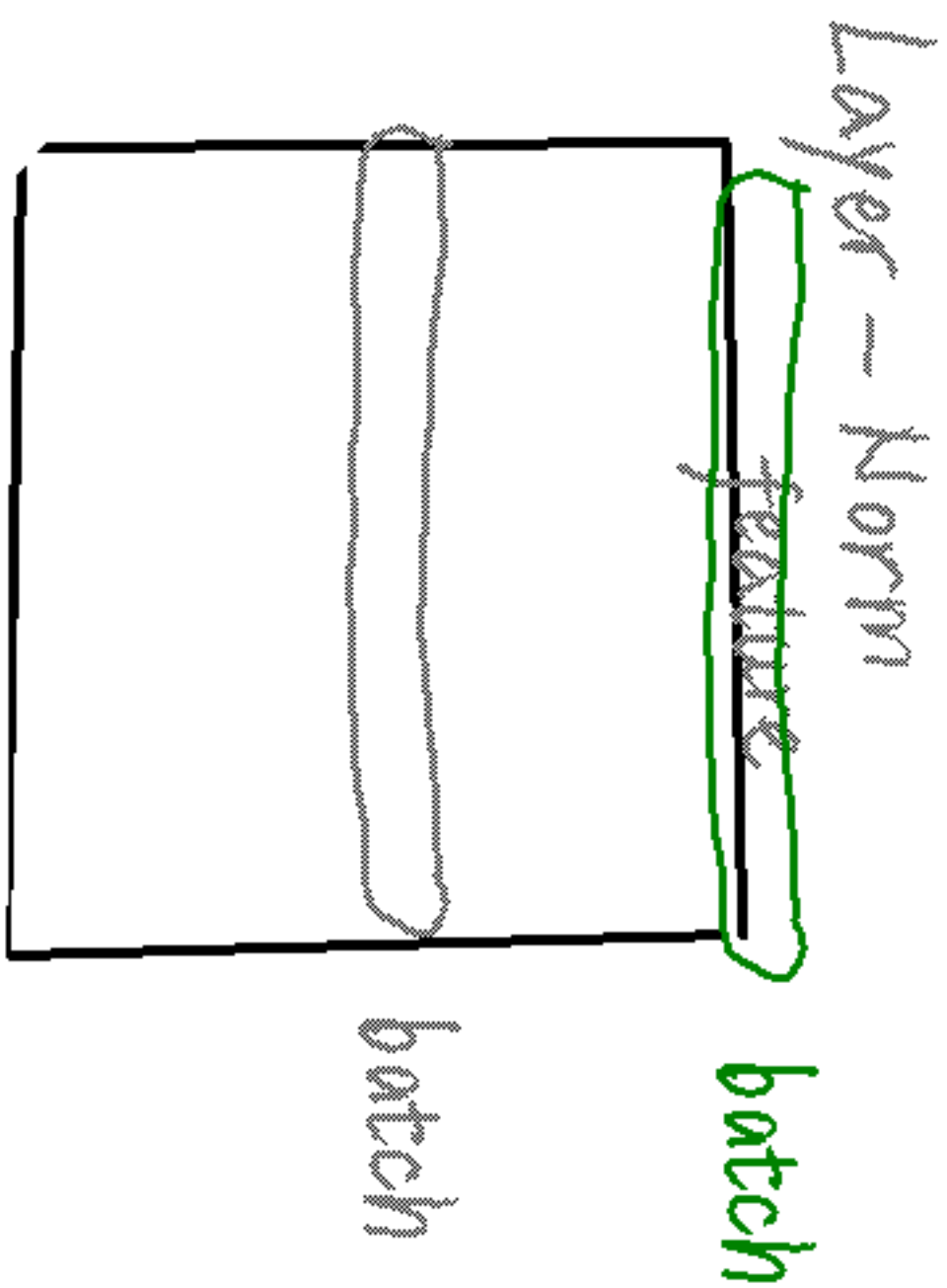


每次把每个列(每个特征)  
在一个 mini-batch 里, 均值变成 0, 方差变为 1  
(训练)

把全局均值和方差计算出 (预测时)

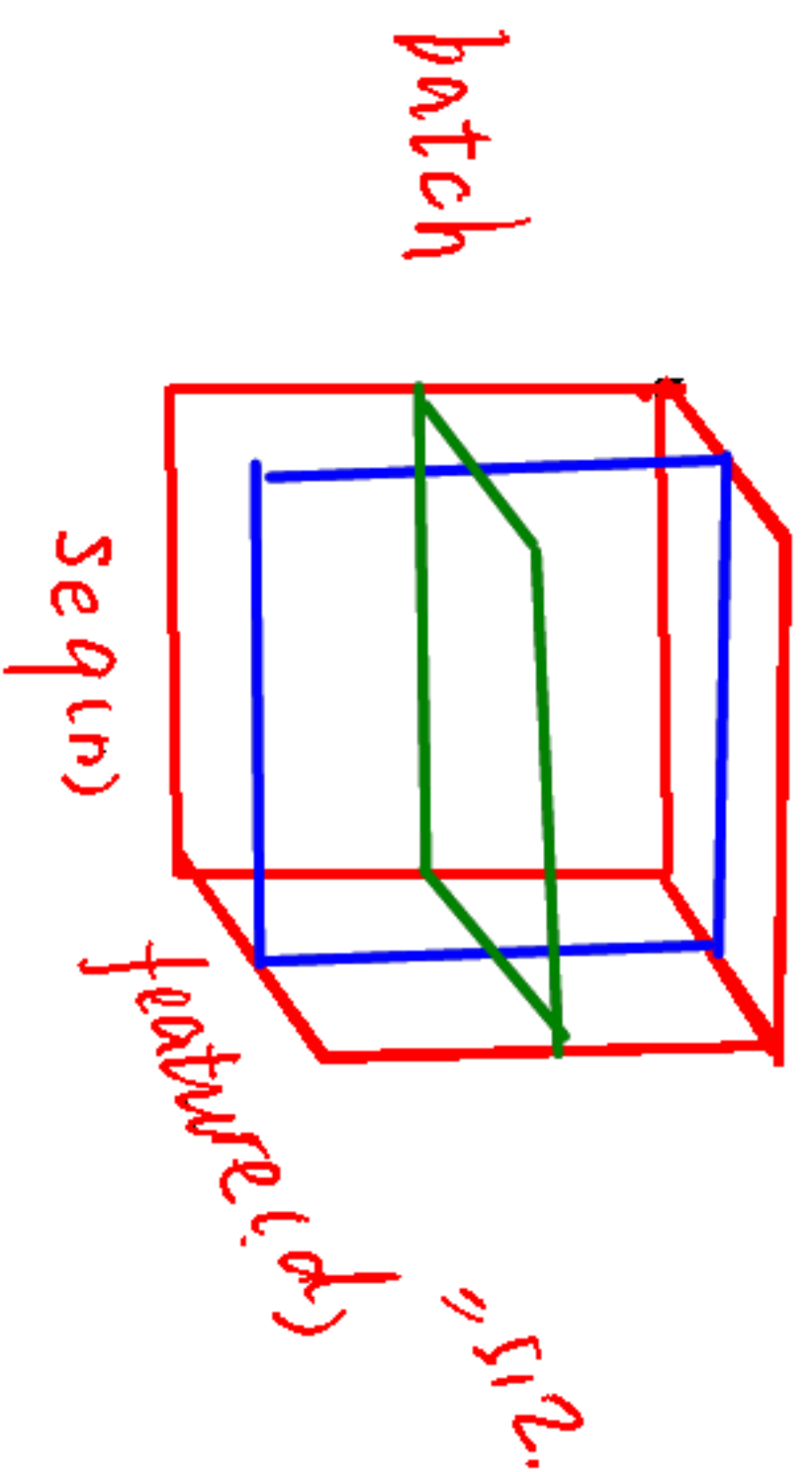
同时学习  $\lambda$  和  $\beta$ , 通过学习可以把向量放成任意均值, 方差的东西.

## Layer-Norm



把每行(每个样本)作 Normalization.

把每行(每个样本)作 Normalization.

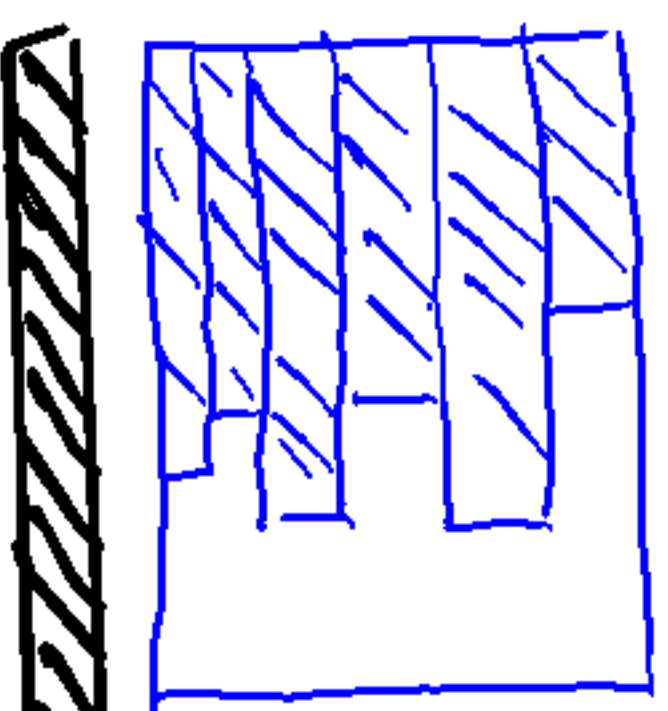


每个词对应一个 seq.

蓝色 = Batch Norm

绿色 = Layer Norm

实际情况, seq 长度可能会变化

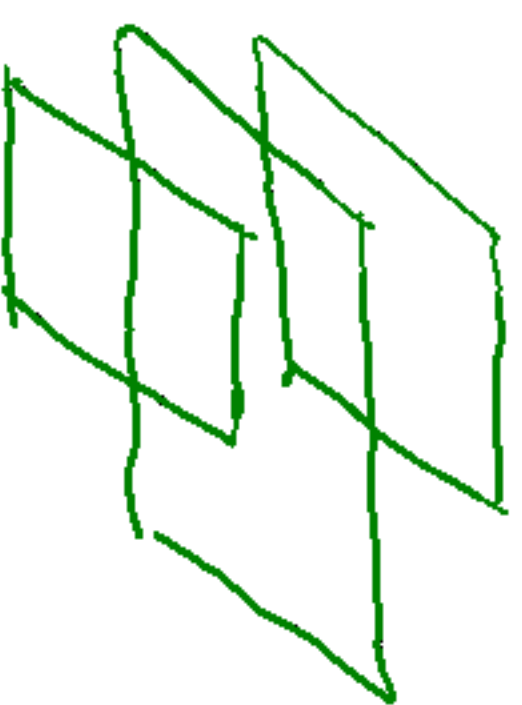


← Batch-Norm



→ 没见过 同时, 预测时用全局的 Mean 和 Var

样本长度变化大时  
Mean, Var 抖动大



← Layer Norm

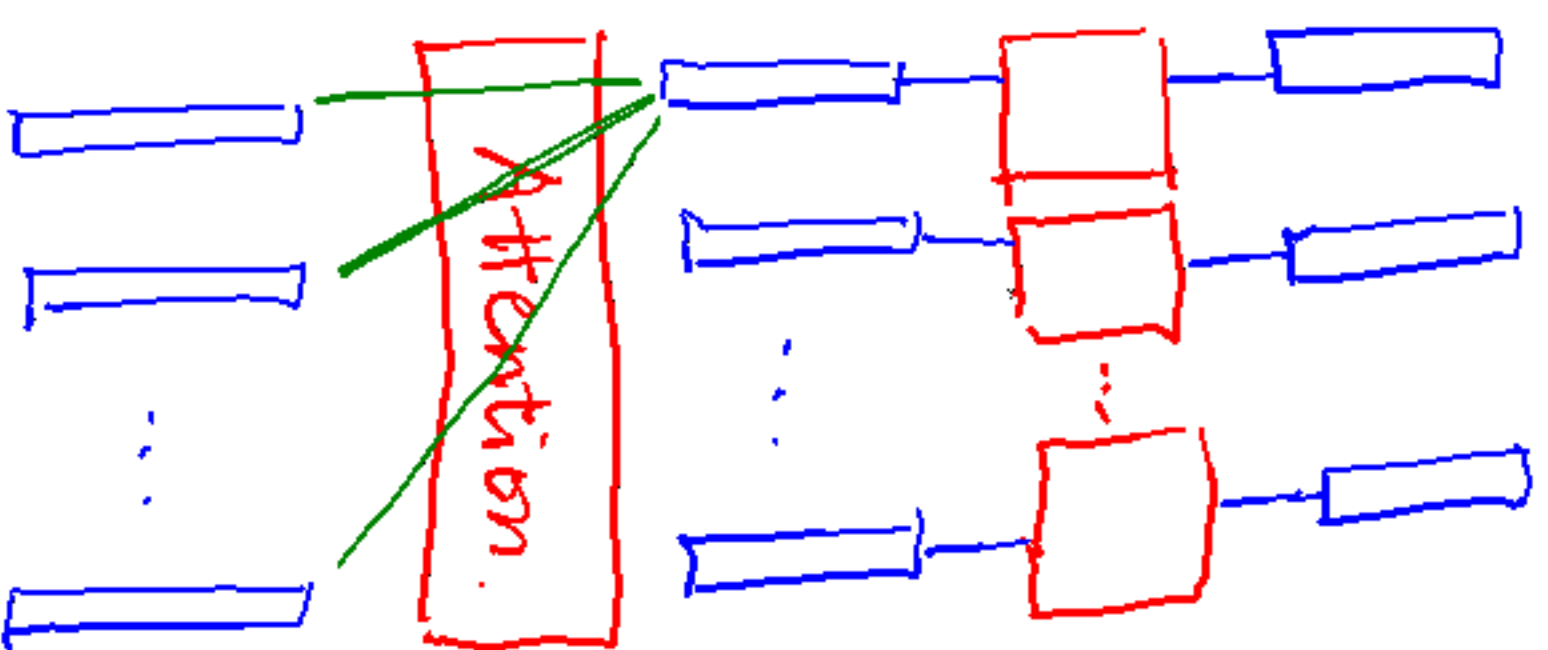
每个样本计算 Mean 和 Var

也不需要全局的 Mean 和 Var

Same: RNN 和 Transformer 都是用 一个 MLP / 线性层  
做 语义空间转换

Difference: 如何传递序列信息

RNN = 上时刻  $\rightarrow$  下时刻  
Transformer = Attention 层, 全局信息  $\rightarrow$  MLP 转换

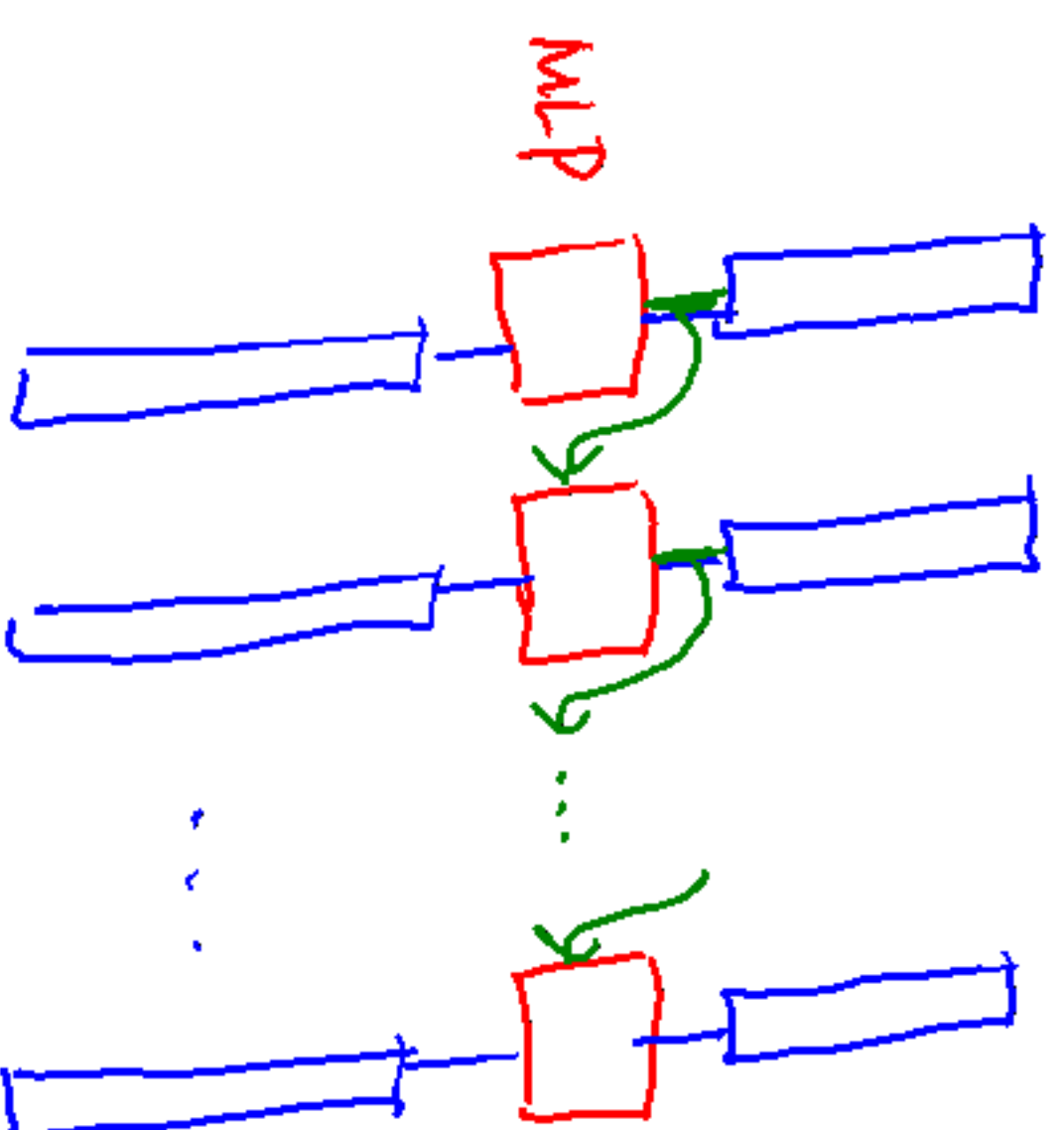


权重一样

MLP - 可分开做

$\therefore$  MLP 可分开

加权 序列信息汇聚完成



权重一样

RNN