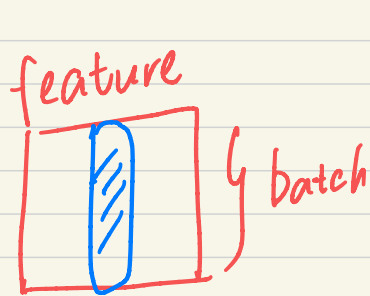
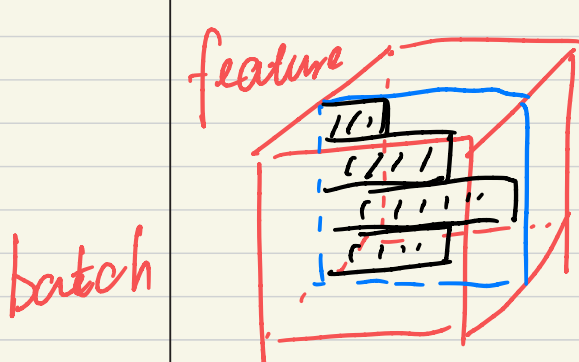
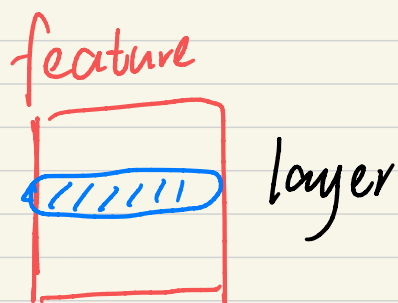


Batch normalization

Layer normalization

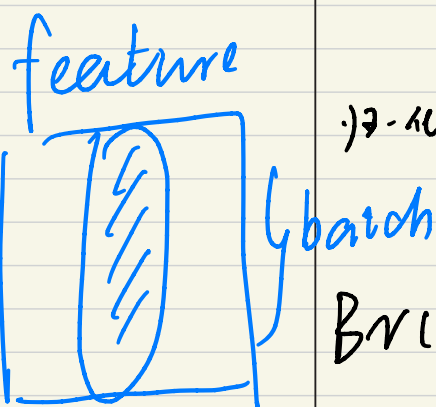
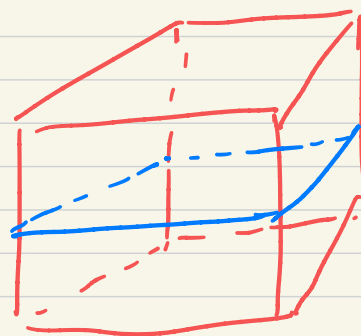


均值为0 方差为1



Seq (n)

Seq 长度不统一



归一化的作用: 特征输入: 激活函数前
避免落在饱和区

$$BN(x_i) = \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}} + \beta \rightarrow \text{满足标准分布}$$

极限: 归一化后方差为1

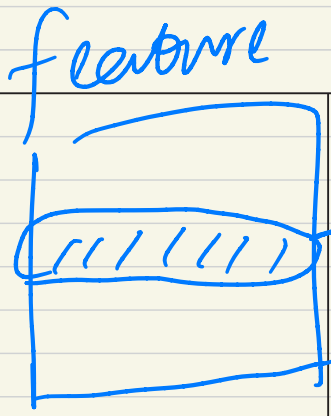
α : 控制归一化后的方差

β : 控制归一化后的均值

$$X \sim N(0, 1)$$

$$Y = \alpha X + \beta \sim N(\beta, \alpha^2)$$

$$L_N(x_i) = \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \rho}} + \beta \rightarrow \text{所有特征}$$



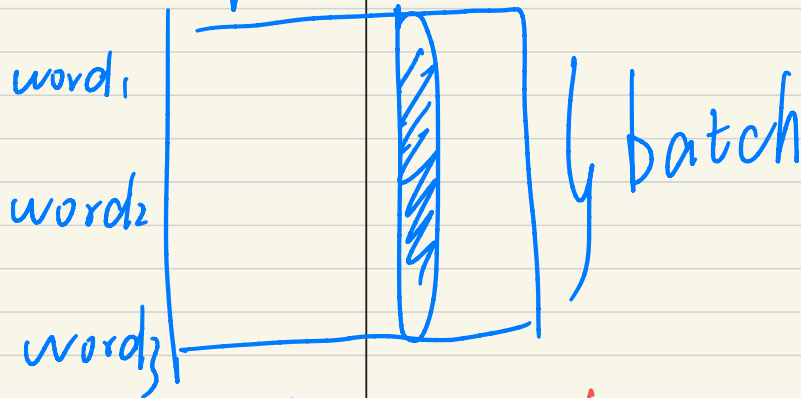
Sample

(适用变长数据层)

二维: 字的情绪划分

feature B_N

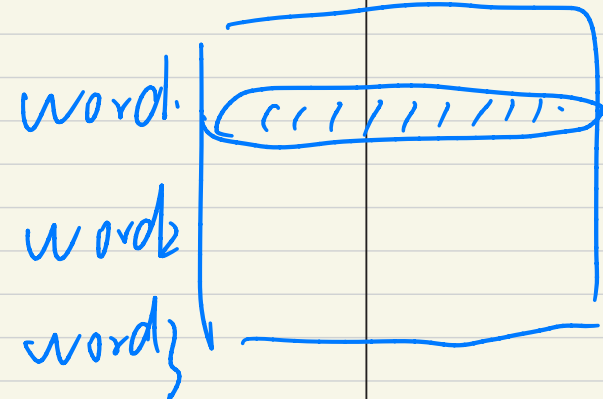
计算每个特征每个批量下的均值和方差



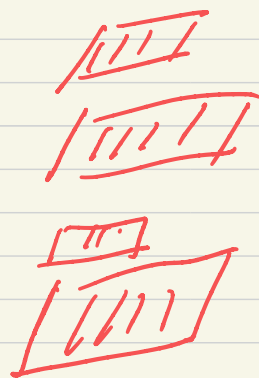
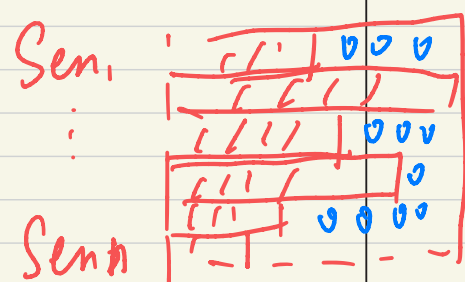
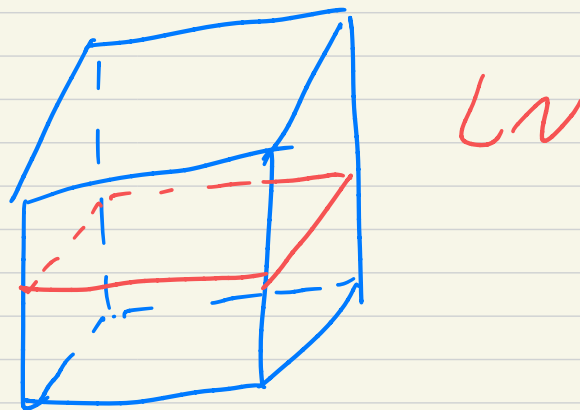
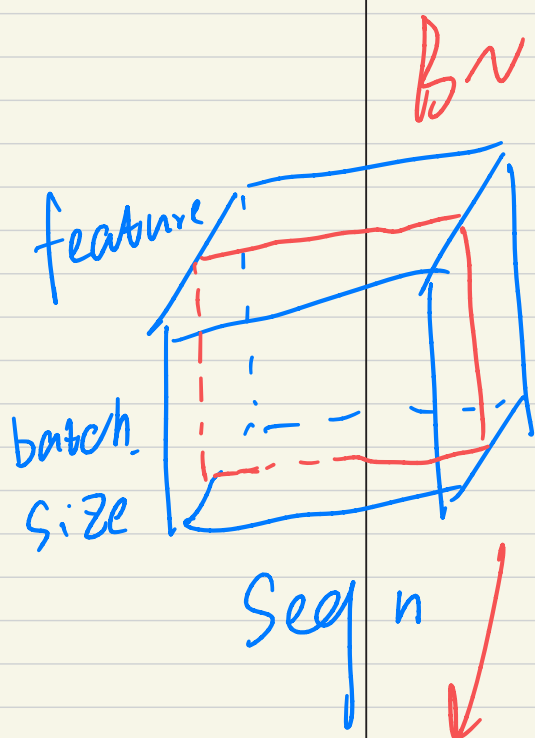
batch

feature L_N

计算每个样本所有特征的均值和方差



三维下: 机器翻译



不同句子的词在某个特征上的取值不一定有

按 batch 数据后针对同特征
的均值和方差波动很大

新句子很长的话不适用原来的均值和方差

不需要存全局的均值和方差
而且只针对该样本更稳定

切片后可以指定归一化的 shape

可以选择是对最后一个维度归一化还是对整片数据归一化

`nn.LayerNorm(embedding-dim)`
/ (5, 10)