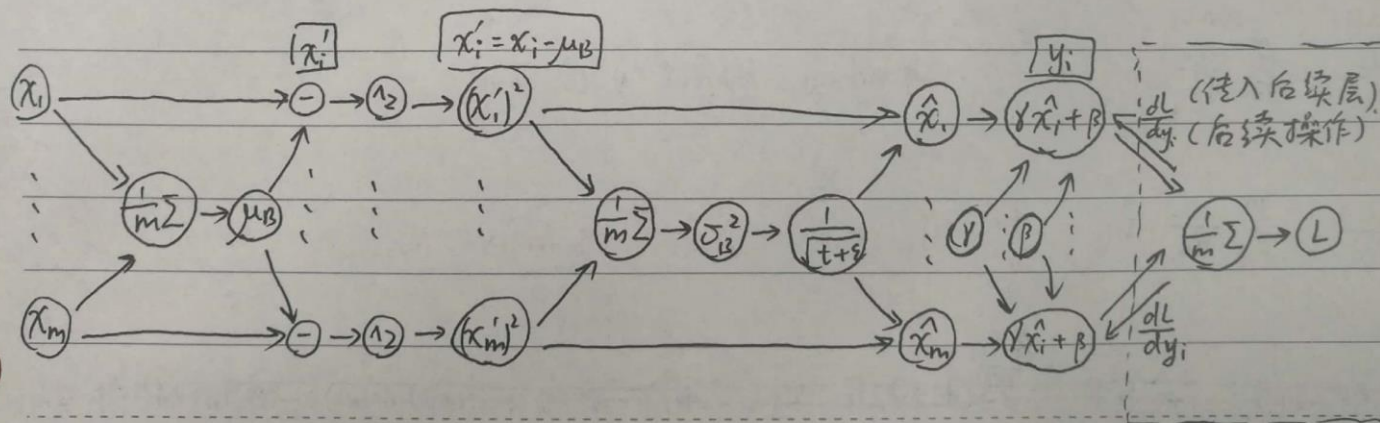
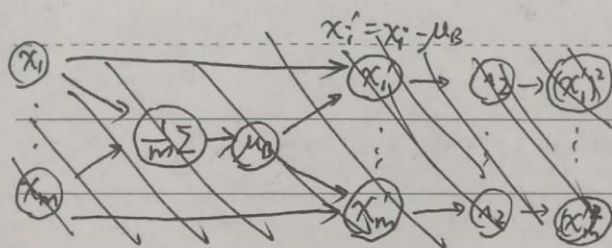


# ① Normalization

- Batch normalization, Layer normalization, Spatial normalization, Group normalization.
- 区别仅在于用在标准化上的维度不同, 写出BN的BP后, 对输入维度进行转置直接调BN包即可实现后三者 (详见代码).

先画出一维  $m$  个样本的 computation graph: ~~U 样率~~



从  $\frac{dL}{dy_i}$  开始, 进行BP, 可参照原论文对比求导正确性. (红笔向量化)

$$\frac{dL}{d\hat{x}_i} = \frac{dL}{dy_i} \cdot \gamma \quad \Rightarrow \text{broadcast 即可}$$

$$\frac{dL}{d\sigma_B^2} = \sum_i \frac{dL}{d\hat{x}_i} \cdot \frac{d\hat{x}_i}{d\sigma_B^2} = \sum_{i=1}^m \frac{dL}{d\hat{x}_i} \left( -\frac{1}{2} \frac{1}{\sigma_B^2 + \epsilon} \cdot \hat{x}_i \right) \quad \Rightarrow \text{括号内 broadcast, } \sum_{i=1}^m \text{用 } (N \times D)^T (N \times D) \text{ 实现.}$$

$$\frac{dL}{d\hat{x}_i'} = \frac{dL}{d\hat{x}_i} \cdot \frac{d\hat{x}_i}{d\hat{x}_i'} + \frac{dL}{d\sigma_B^2} \cdot \frac{d\sigma_B^2}{d\hat{x}_i'} = \frac{dL}{d\hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{dL}{d\sigma_B^2} \cdot \frac{1}{m} 2\hat{x}_i' \quad \Rightarrow \text{broadcast 即可.}$$

$$\frac{dL}{du_B} = \sum_i \frac{dL}{d\hat{x}_i'} \cdot \frac{d\hat{x}_i'}{du_B} = \sum_{i=1}^m \frac{dL}{d\hat{x}_i'} \cdot (-1) \quad \Rightarrow \sum_{i=1}^m \text{用 } \text{np.ones}(N) \cdot \text{dot}(N \times D) \text{ 实现.}$$

$$\frac{dL}{d\pi_i} = \frac{dL}{du_B} \cdot \frac{du_B}{d\pi_i} + \frac{dL}{d\hat{x}_i'} \cdot \frac{d\hat{x}_i'}{d\pi_i} = \frac{1}{m} \frac{dL}{du_B} + \frac{dL}{d\hat{x}_i'} \quad \Rightarrow \text{broadcast 即可.}$$

$$\frac{dL}{d\beta} = \sum_{i=1}^m \frac{dL}{dy_i} \quad \Rightarrow \sum_{i=1}^m \text{用 } \text{np.ones}(N) \cdot \text{dot}(N \times D) \text{ 实现}$$

$$\frac{dL}{d\gamma} = \sum_{i=1}^m \frac{dL}{dy_i} \cdot \hat{x}_i \quad \Rightarrow \sum_{i=1}^m \text{用 } (N \times D)^T (N \times D) \text{ 实现, 最后取对角}$$

关于取对角:  $\sigma_B^2$  与  $\gamma$  为特征独有, 相互之间不干扰,  $(N \times D)^T (N \times D)$  得出的种角对角

项代表

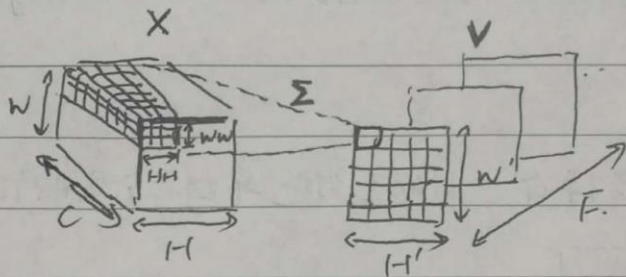
关于取对角: 等价于两矩阵作 element-wise 的乘法, 再用向量化实现  $\sum_{i=1}^m$ .

$$\text{即 } \frac{dL}{d\gamma} = \text{np.ones}(N) \cdot \text{dot} \left( \frac{dL}{dy} * \hat{x} \right) \\ (= \text{np.ones}(N) \cdot \text{dot} \left( \frac{dL}{dy} \odot \hat{x} \right))$$

$$\frac{dL}{d\sigma_B^2} \text{ 同理}$$



## ② CNN



前向传导简单，反向传导需理清各变量依赖关系。

• 参数：卷积核  $W_F \in \mathbb{R}^{H' \times W' \times C}$  可见其与  $F$  有关，所以先讨论  $V, F$ 。

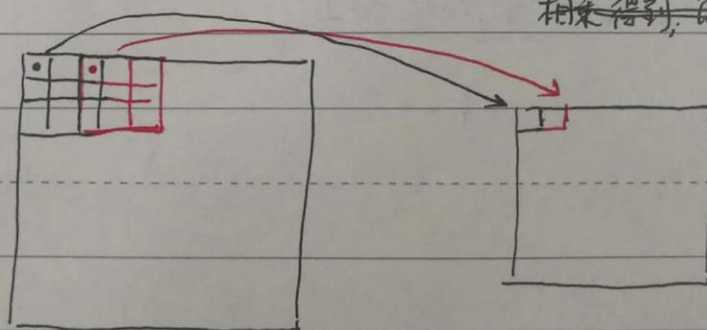
但  $W_F \in \mathbb{R}^{H' \times W' \times C}$ ，注意到  $C$  维上行为是相同的，所以再简化，讨论  $V, C$ 。

故  $V, F, C$ ，(1)  $\frac{dL}{dw_{hw}} = \sum_{k,l} \frac{dL}{dv_{kl}} \cdot \frac{dv_{kl}}{dw_{hw}}$  (下标与  $w$  重了... 不多用 orz)

$$= \sum_l \sum_k \left| \frac{dL}{dv_{kl}} \right| \cdot \left| \frac{dv_{kl}}{dw_{hw}} \right|$$

↑ 层传入

找出  $V$  中与  $w_{hw}$  有关的元素，(而  $V$  与  $X$  相乘得到，即寻找  $V$  中  $X$  中与  $w_{hw}$  相乘过的元素)



由于  $V = X * W$ ，即在  $X$  中寻找与  $w_{hw}$  相乘过的元素

$$\text{因此 } \frac{dL}{dw} = X_s * \frac{dL}{dw}$$

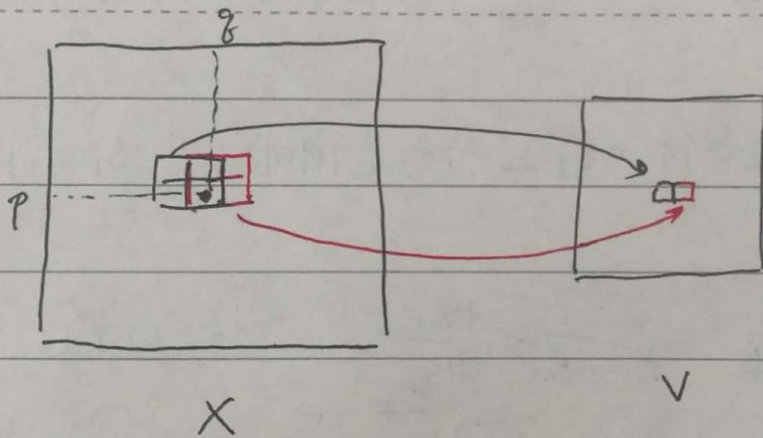
( $*$  为卷积操作)

其中  $X_s$  是按步长间隔取出的矩阵，但卷积步长为 1

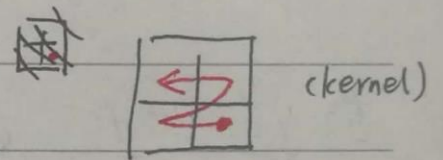
(2)  $\frac{dL}{db} = \sum_{k,l} \frac{dL}{dV_{kl}} \cdot \frac{dV_{kl}}{db}$  , 直接将  $\frac{dL}{dV}$  所有元素求和即可.

(3) 求  $\frac{dL}{dx}$  . 由于涉及多样本, 此时先讨论单样本, 再进行向量化

$$\frac{dL}{dx_{pg}^{(i)}} = \sum_{k,l} \underbrace{\frac{dL}{dV_{kl}}}_{\text{上层传入}} \underbrace{\frac{dV_{kl}}{dx_{pg}^{(i)}}}_{\text{找出V中与 } x_{pg}^{(i)} \text{ 有关的元素}}$$



可见, 与  $x_{pg}^{(i)}$  相关的  $w$  顺序是逆的, 即



(“逆”是相对卷积得出  
的值在  $V$  中的顺序而言)

因此  $\frac{dL}{dx} = \text{rot}180(w) * \frac{dL}{dV}$

其中卷积步长为1

#### (4) 边界情况

- 对  $\frac{dL}{dW}$  无 0 填充. (但注意  $\frac{dL}{dW} = X_s * \frac{dL}{dV}$  中,  $X_s$  从填充后的  $X$  中取出, 即前向传播时用的  $X$ ).

- 对  $\frac{dL}{dX}$  需在  $\frac{dL}{dV}$  边缘及内部 0 填充.

边缘填充 "W 的长度 - 1" (或 "宽度 - 1"), 内部填充 "步长 - 1"

同时注意, 求出的  $\frac{dL}{dX}$  是带填充的, 须去掉边缘

填充细节是通过分析  $X * W = V$  第一行的行为总结出的, 很有意思, 留作思考. 方法类似 (11)、(13) 的画图分析.

- 这里给出  $\frac{dL}{dV}$  填充后大小与  $\text{rot180}(W) * \frac{dL}{dV}$  的大小

设  $X \in \mathbb{R}^{H \times W}$ ,  $W \in \mathbb{R}^{H_H \times W_W}$ ,  $\frac{dL}{dV} \in \mathbb{R}^{H' \times W'}$

步长为  $S$ , 前导时对  $X$  填充系数为  $P$ .

则填充后  $X_{\text{-pad}}$  长为:  ~~$(H+2P)/S+1$~~   $H+2P$ .

按上述规则填充后  $\frac{dL}{dV}_{\text{-pad}}$  长为:  $(H'-1)(S-1) + H' + 2(CHH-1)$

$\text{rot180}(W) * \frac{dL}{dV}$  后 (注意步长始终为 1) 长为:

$$(H'-1)S + 1 + 2(CHH-1) - HH + 1 = (H'-1)S + HH$$

~~注意到  $W * X_{\text{-pad}}$  的长为  $H+2P$~~

注意到  $W * X_{\text{-pad}}$  的长为:  $(H+2P - HH)/S + 1 = H'$ , 故

$$\rightarrow H+2P = (H'-1)S + HH$$

可见  $\frac{dL}{dV}$  填充后再卷积, 大小与  $X_{\text{-pad}}$  相同, 这也说明反向传播的维度上是正确/合理的. 前述的.