

量化参数选择

Sensetime HPC Group

定义

- 为了研究方便，我们假定： $x_i \sim N(0,1)$ 且独立同分布

- 我们定义量化函数 $Q(x, s) = \begin{cases} [\frac{x}{s}], & \text{else} \\ 127, & \frac{x}{s} > 127.5 \\ -127, & \frac{x}{s} < -127.5 \end{cases}$, 额外定义截断值 $c = s * 127.5$

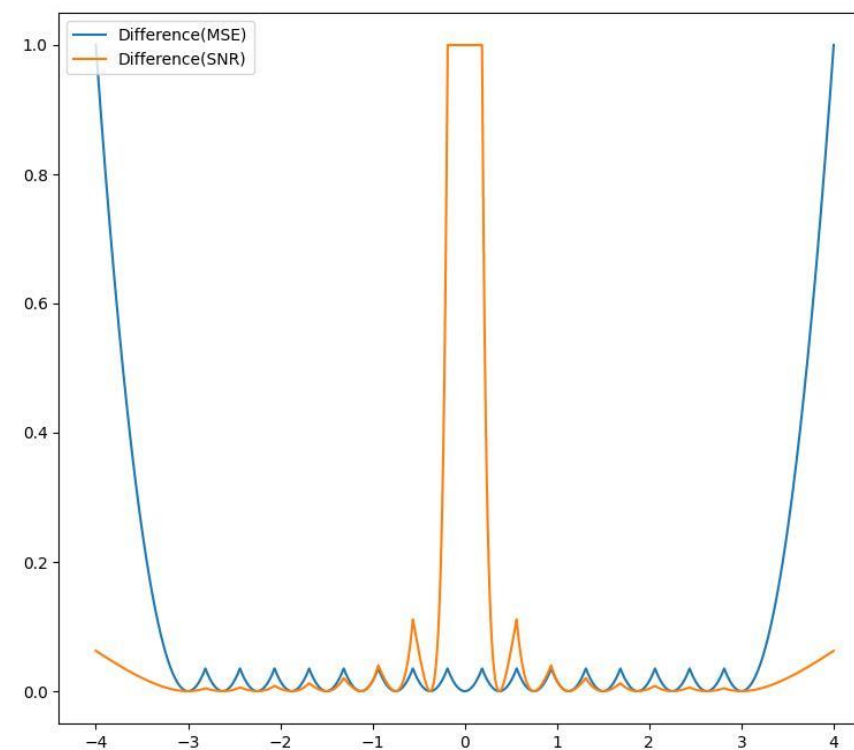
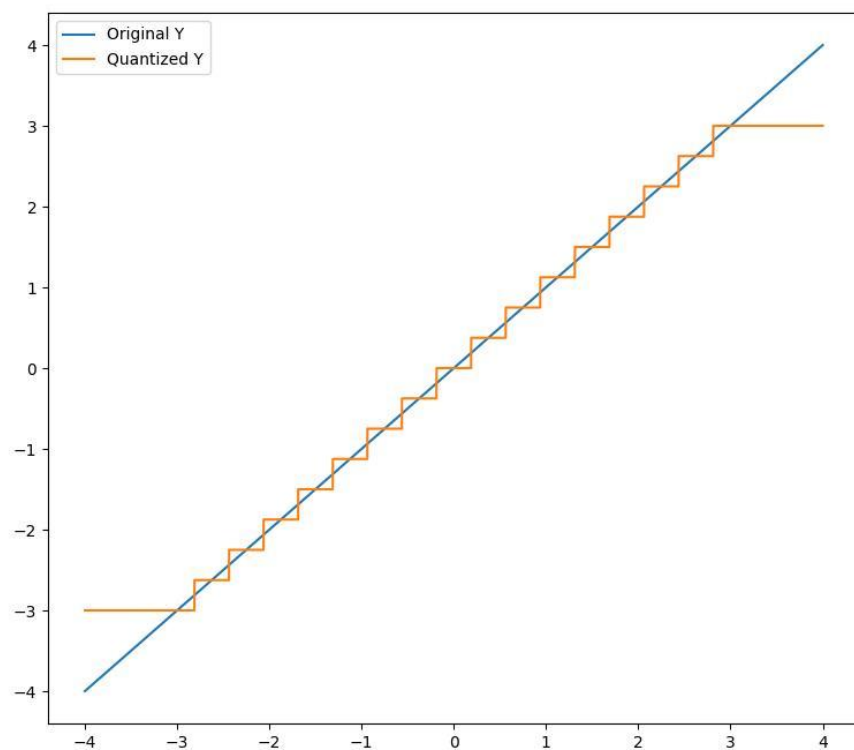
- 我们定义反量化函数 $DQ(x, s) = x * s$

- 我们定义x的量化值为 $x'(s) \triangleq DQ(Q(x, s), s)$

- $L_{MSE} = \frac{\sum_i (x_i - x_i')^2}{N}$

- $L_{N/S} = \frac{\sum_i \frac{(x_i - x_i')^2}{x_i^2}}{N}$

量化误差分析 - 量化函数与量化误差



量化计算例子

$$\bullet \begin{bmatrix} 1.3 & 4.7 & -0.5 \\ 2.1 & 6.0 & -1.1 \\ 10.0 & 0.3 & 25.1 \end{bmatrix} \xrightarrow{Q(x,0.1)} \begin{bmatrix} 13 & 47 & -5 \\ 21 & 6 & -11 \\ 100 & 3 & 127 \end{bmatrix} \xrightarrow{DQ(x,0.1)} \begin{bmatrix} 1.3 & 4.7 & -0.5 \\ 2.1 & 6.0 & -1.1 \\ 10.0 & 0.3 & 12.7 \end{bmatrix}$$

$$\bullet L_{MSE} = \frac{\sum_i (x_i - x_{i'})^2}{N} = \frac{\left\| \begin{bmatrix} 1.3 & 4.7 & -0.5 \\ 2.1 & 6.0 & -1.1 \\ 10.0 & 0.3 & 25.1 \end{bmatrix} - \begin{bmatrix} 1.3 & 4.7 & -0.5 \\ 2.1 & 6.0 & -1.1 \\ 10.0 & 0.3 & 12.7 \end{bmatrix} \right\|_2^2}{9} = \frac{12.4^2}{9} = 17.08$$

$$\bullet L_{N/S} = \frac{\sum_i \frac{(x_i - x_{i'})^2}{x_i^2}}{N} = \frac{12.4^2}{9 * 25.1^2} = 0.027$$

量化参数选择

- 量化参数选择：选择一个合适的 s ，使得量化误差最小

$$s^* = \operatorname{argmin}_s \frac{\sum_i (x_i - x_i'(s))^2}{N}$$

- 注意到 s 与 截断值 c 之间存在固定比例关系，因此下文中我们有时也变相去求解最优截断值。

量化参数选择

- 一个简单的思路：

$$L_{MSE} = \frac{\sum_i (x_i - x_i')^2}{N} = \frac{\left\| \begin{bmatrix} 1.3 & 4.7 & -0.5 \\ 2.1 & 6.0 & -1.1 \\ 10.0 & 0.3 & 25.1 \end{bmatrix} - \begin{bmatrix} 1.3 & 4.7 & -0.5 \\ 2.1 & 6.0 & -1.1 \\ 10.0 & 0.3 & 12.7 \end{bmatrix} \right\|_2^2}{9} = \frac{12.4^2}{9} = 17.08$$

- MSE误差主要来自于网络中较大的项，这是因为边界值的MSE是发散的。因此我们可以设计足够大的s，使得所有边界值都可以被表示：

$$s_{MAX} = \frac{Max(Abs(X))}{127} = \frac{25.1}{127} = 0.2$$

量化参数选择

$$\bullet \begin{bmatrix} 1.3 & 4.7 & -0.5 \\ 2.1 & 6.0 & -1.1 \\ 10.0 & 0.3 & 25.1 \end{bmatrix} \xrightarrow{Q(x,0.1)} \begin{bmatrix} 13 & 47 & -5 \\ 21 & 60 & -11 \\ 100 & 3 & 127 \end{bmatrix} \xrightarrow{DQ(x,0.1)} \begin{bmatrix} 1.3 & 4.7 & -0.5 \\ 2.1 & 6.0 & -1.1 \\ 10.0 & 0.3 & 12.7 \end{bmatrix}$$

$$\bullet \begin{bmatrix} 1.3 & 4.7 & -0.5 \\ 2.1 & 6.0 & -1.1 \\ 10.0 & 0.3 & 25.1 \end{bmatrix} \xrightarrow{Q(x,0.2)} \begin{bmatrix} 6 & 24 & -2 \\ 10 & 30 & -6 \\ 50 & 2 & 126 \end{bmatrix} \xrightarrow{DQ(x,0.2)} \begin{bmatrix} 1.2 & 4.8 & -0.4 \\ 2.0 & 6.0 & -1.2 \\ 10.0 & 0.4 & 25.2 \end{bmatrix}$$

$$\bullet L_{MSE} = \frac{\sum_i (x_i - x_{i'})^2}{N} = \frac{\left\| \begin{bmatrix} 1.3 & 4.7 & -0.5 \\ 2.1 & 6.0 & -1.1 \\ 10.0 & 0.3 & 25.1 \end{bmatrix} - \begin{bmatrix} 1.2 & 4.8 & -0.4 \\ 2.0 & 6.0 & -1.2 \\ 10.0 & 0.4 & 25.2 \end{bmatrix} \right\|_2^2}{9} = \frac{7 * 0.1^2}{9} = 0.078$$

$$\bullet L_{N/S} = \frac{\sum_i \frac{(x_i - x_{i'})^2}{x_i^2}}{N} = ?$$

量化参数选择 - 最大值截断

- 方案1：最大值截断：

$$s_{MAX} = \frac{Max(Abs(X))}{127} = \frac{25.1}{127} = 0.2$$

- 问题：在已知 $x_i \sim N(0,1)$ 且独立同分布的情况下，假设总体 X 元素数量趋近于无穷大，试问 $E\{Max(X)\}$ 是否同样趋近于无穷大？在此情况下 s_{MAX} 是否也同样趋近于无穷大？

$$\begin{aligned} E\{Max(X)\} &= \int_{-inf}^{+inf} u P(Max(X) = u) du \\ &= C_N^1 \int_{-inf}^{+inf} u \varphi(u) \left(\int_{-inf}^u \varphi(v) dv \right)^{N-1} du \end{aligned}$$

量化参数选择 - 最大值截断

$$\lim_{N \rightarrow \infty} (C_N^1 \int_{-\infty}^{+\infty} u \varphi(u) (\int_{-\infty}^u \varphi(v) dv)^{N-1} du) = \sqrt{2 \log(N)} = \infty$$

$$s_{MAX} = \frac{\sqrt{2 \log(N)}}{127} = \frac{\infty}{127} = \infty$$

will have: $\lim_{N \rightarrow \infty} (E\{L_{MSE}\}) = \infty$

- 也就是说最大值截断在元素数量趋于无限时，会出现误差发散的情况。

量化参数选择 - 分位数截断

- 方案2: $k - \sigma$ 截断:

$$s_{k-\sigma} = \frac{k * \sigma}{127}$$

意义: σ 与 N 无关, 因此当 N 趋近于无穷大时, 量化误差的期望是收敛的。

- 证明:

$$E\{L_{MSE}\} = \int_{-k\sigma - \frac{s}{2}}^{k\sigma + \frac{s}{2}} (x' - x)^2 \varphi(x) dx + \int_{k\sigma + \frac{s}{2}}^{inf} (x - k\sigma)^2 \varphi(x) dx + \int_{-inf}^{-k\sigma - \frac{s}{2}} (x + k\sigma)^2 \varphi(x) dx$$

$$E\{L_{MSE}\} = \int_{-k\sigma - \frac{s}{2}}^{k\sigma + \frac{s}{2}} (x' - x)^2 \varphi(x) dx + \int_{k\sigma + \frac{s}{2}}^{inf} (x - k\sigma)^2 \varphi(x) dx + \int_{-inf}^{-k\sigma - \frac{s}{2}} (x + k\sigma)^2 \varphi(x) dx$$

- 上式中第一项称为表示误差，显然存在上界，后两项称为正负截断误差，只需要证明后两项存在上界。

$$\int_{-inf}^{-k\sigma - \frac{s}{2}} (x + k\sigma)^2 \varphi(x) dx = \int_{k\sigma + \frac{s}{2}}^{inf} (x - k\sigma)^2 \varphi(x) dx$$

- 注意到在高斯分布中，正负截断误差相等。

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$$\int_{k\sigma + \frac{s}{2}}^{inf} (x - k\sigma)^2 \varphi(x) dx = \frac{1}{\sqrt{2\pi}} \int_{k\sigma + \frac{s}{2}}^{inf} (x - k\sigma)^2 e^{-\frac{x^2}{2}} dx$$

$$\frac{1}{\sqrt{2\pi}} \int_{k\sigma + \frac{s}{2}}^{inf} (x - k\sigma)^2 e^{-\frac{x^2}{2}} dx = \int_{k\sigma + \frac{s}{2}}^{inf} (x^2 - 2k\sigma x + k^2\sigma^2) e^{-\frac{x^2}{2}} dx$$

$$= \int_{k\sigma + \frac{s}{2}}^{inf} x^2 e^{-\frac{x^2}{2}} dx - 2k\sigma \int_{k\sigma + \frac{s}{2}}^{inf} x e^{-\frac{x^2}{2}} dx + k^2\sigma^2 \int_{k\sigma + \frac{s}{2}}^{inf} e^{-\frac{x^2}{2}} dx$$

$$= \int_{k\sigma + \frac{s}{2}}^{\inf} x^2 e^{-\frac{x^2}{2}} dx - 2k\sigma \int_{k\sigma + \frac{s}{2}}^{\inf} x e^{-\frac{x^2}{2}} dx + k^2\sigma^2 \int_{k\sigma + \frac{s}{2}}^{\inf} e^{-\frac{x^2}{2}} dx$$

$$x = \sqrt{2}t, dx = \sqrt{2}dt$$

$$\int_{\frac{\sqrt{2}}{2}(k\sigma + \frac{s}{2})}^{\inf} t^2 e^{-t^2} dt - 2k\sigma \int_{\frac{\sqrt{2}}{2}(k\sigma + \frac{s}{2})}^{\inf} t e^{-t^2} dt + k^2\sigma^2 \int_{\frac{\sqrt{2}}{2}(k\sigma + \frac{s}{2})}^{\inf} e^{-t^2} dt$$

$$\int_{\frac{\sqrt{2}}{2}(k\sigma + \frac{s}{2})}^{inf} t^2 e^{-t^2} dt - 2k\sigma \int_{\frac{\sqrt{2}}{2}(k\sigma + \frac{s}{2})}^{inf} t e^{-t^2} dt + k^2\sigma^2 \int_{\frac{\sqrt{2}}{2}(k\sigma + \frac{s}{2})}^{inf} e^{-t^2} dt$$

$$\int_0^x \frac{2}{\sqrt{\pi}} e^{-t^2} dt = erf(x)$$

$$\sqrt{\frac{\pi}{2}} (k^2\sigma^2 + 1) erf\left(\frac{x}{\sqrt{2}}\right) - (x - 2k\sigma) e^{-\frac{x^2}{2}} \Bigg|_{k\sigma + \frac{s}{2}}^{inf}$$

$$\sqrt{\frac{\pi}{2}} (k^2 \sigma^2 + 1) \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right) - (x - 2k\sigma) e^{-\frac{x^2}{2}} \Big|_{k\sigma + \frac{s}{2}}^{inf}$$

$$x = inf \rightarrow \sqrt{\frac{\pi}{2}} (k^2 \sigma^2 + 1)$$

$$x = k\sigma + \frac{s}{2} \rightarrow \sqrt{\frac{\pi}{2}} (k^2 \sigma^2 + 1) \operatorname{erf}\left(\frac{k\sigma + \frac{s}{2}}{\sqrt{2}}\right) + (k\sigma - \frac{s}{2}) e^{-\frac{(k\sigma + \frac{s}{2})^2}{2}}$$

$$\sqrt{\frac{\pi}{2}} (k^2 \sigma^2 + 1) \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right) - (x - 2k\sigma) e^{-\frac{x^2}{2}} \Big|_{k\sigma + \frac{s}{2}}^{inf} < \sqrt{\frac{\pi}{2}} (k^2 \sigma^2 + 1)$$

截断误差上界

$$E\{L_{MSE}\} = \underbrace{\int_{-k\sigma - \frac{s}{2}}^{k\sigma + \frac{s}{2}} (x' - x)^2 \varphi(x) dx}_{\text{有界}} + \underbrace{\int_{k\sigma + \frac{s}{2}}^{inf} (x - k\sigma)^2 \varphi(x) dx}_{\text{有界}} + \underbrace{\int_{-inf}^{-k\sigma - \frac{s}{2}} (x + k\sigma)^2 \varphi(x) dx}_{\text{有界}}$$

有界

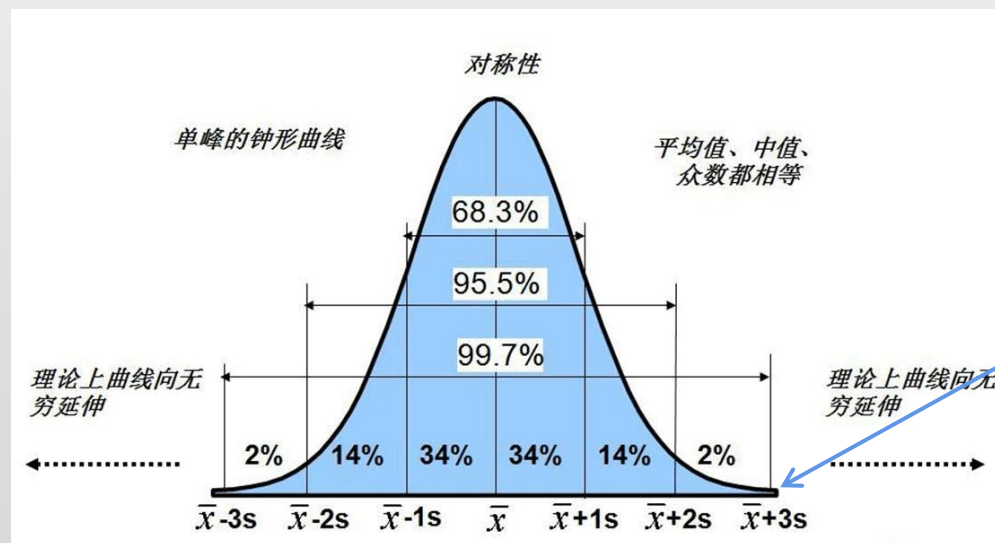
量化参数选择 - 分位数截断

- 方案2: $k - \sigma$ 截断:

$$s_{k-\sigma} = \frac{k * \sigma}{127}$$

好处: σ 与 N 无关, 因此当 N 趋近于无穷大时, 量化误差的期望是有界的。(可以进一步证明收敛性)

在实践中, 我们往往不使用 $k\sigma$ 而是利用分位点与 $k\sigma$ 之间的关系确定截断值。



以第99.99%大的值作为截断值

这个值接近于 4σ , 此时的 $s = \frac{4\sigma}{127}$

量化参数选择 - 最优截断

- 方案3：最优截断：

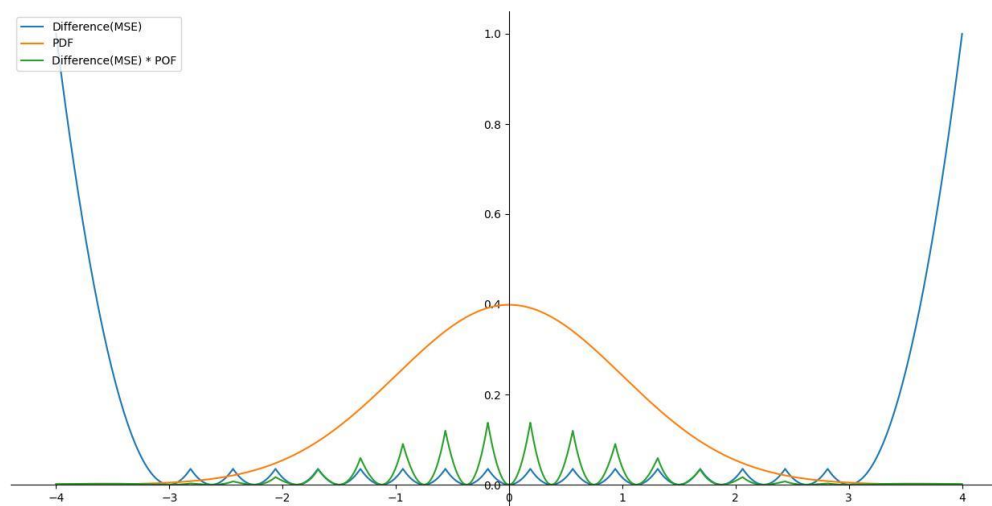
$$s_{opt} = \underset{s}{\operatorname{argmin}} E\{L_{MSE}\}$$

$$E\{L_{MSE}\} = \int_{c-\frac{s}{2}}^{c+\frac{s}{2}} (x' - x)^2 p(x) dx + \int_{c+\frac{s}{2}}^{inf} (x - c)^2 p(x) dx + \int_{-inf}^{-c-\frac{s}{2}} (x + c)^2 p(x) dx$$

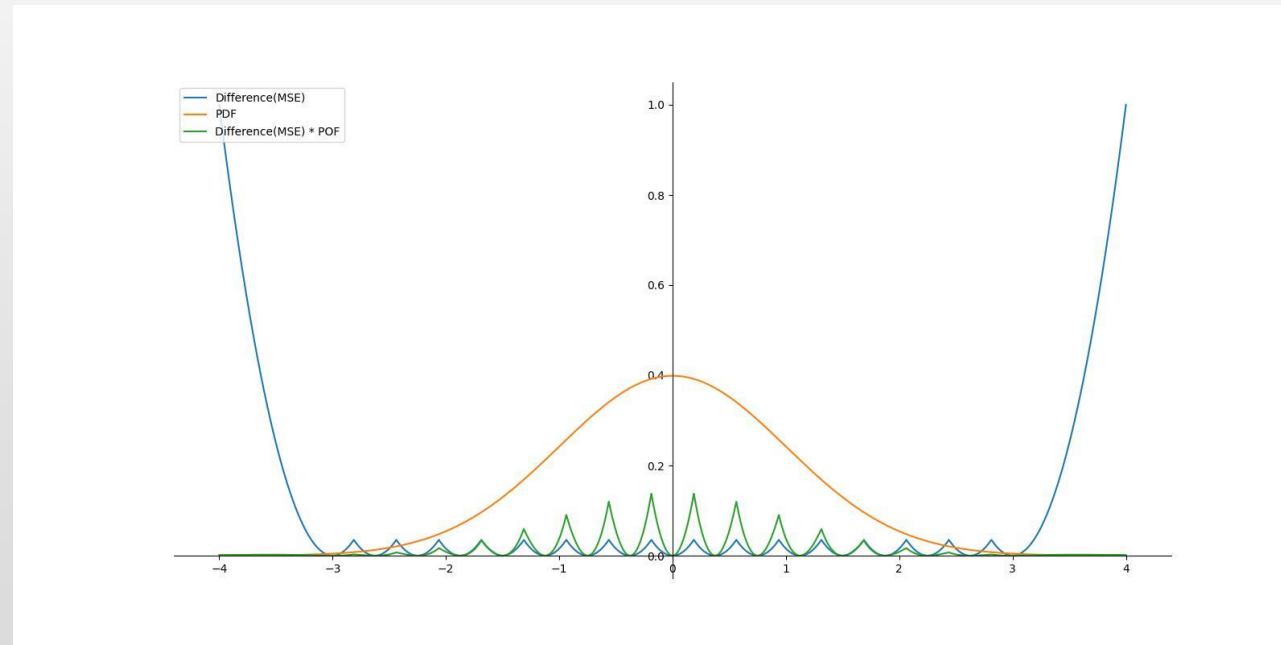
- 思考：我们已经求出了 $E\{L_{MSE}\}$ 的表达式，能否直接解出最优的截断值 c 与尺度因子 s ？
- 此处我们不再局限于高斯分布，而直接求取广义分布 $p(x)$ 的最优截断。

- 首先处理表示误差：

$$\int_{-c-\frac{s}{2}}^{c+\frac{s}{2}} (x' - x)^2 p(x) dx$$



$$\int_{-c-\frac{s}{2}}^{c+\frac{s}{2}} (x' - x)^2 p(x) dx = \sum_i \int_{l_i}^{u_i} (x'_i - x)^2 p(x) dx$$



$$\int_{-c-\frac{s}{2}}^{c+\frac{s}{2}} (x' - x)^2 p(x) dx = \sum_i \int_{l_i}^{u_i} (x'_i - x)^2 p(x) dx$$

- 考虑到 $\varphi(x)$ 原函数不存在，考虑使用常数 $c(i)$ 替换 $\varphi(x)$

$$\sum_i \int_{l_i}^{u_i} (x'_i - x)^2 \varphi(x) dx \approx \sum_i \int_{l_i}^{u_i} (x'_i - x)^2 c(i) dx = \sum_i c(i) \frac{(x - x'_i)^3}{3} \Big|_{l_i}^{u_i}$$

$$\sum_i c(i) \frac{(x - x'_i)^3}{3} \Big|_{l_i}^{u_i} = \frac{s^3}{12} \sum_i c(i) = \frac{s^3}{12} * \frac{P(c + \frac{s}{2}) - P(-c - \frac{s}{2})}{s}$$

$$= P(-c - \frac{s}{2} < x < c + \frac{s}{2}) \frac{s^2}{12}$$

$$\sum_i c(i) = \frac{2P(c + \frac{s}{2}) - P(-c - \frac{s}{2})}{s}$$

- 思考：为什么上式成立？

- 处理截断误差：

$$\int_{c+\frac{s}{2}}^{inf} (x-c)^2 p(x) dx + \int_{-inf}^{-c-\frac{s}{2}} (x+c)^2 p(x) dx$$

$$= \int_{c+\frac{s}{2}}^{inf} (x-c)^2 dP_1(x) = (x-c)^2 P(x) - \int_{c+\frac{s}{2}}^{inf} P_1(x) d(x-c)^2$$

$$= (x-c)^2 P_1(x) - 2((x-c)P_2(x) - P_3(x)) \Big|_{c+\frac{s}{2}}^{inf}$$

$$E\{L_{MSE}\} =$$

$$\begin{aligned} & (x - c)^2 P_1(x) - 2((x - c)P_2(x) - P_3(x)) \Big|_{c + \frac{s}{2}}^{inf} + (x + c)^2 P_1(x) - 2((x + c)P_2(x) - P_3(x)) \Big|_{-inf}^{-c - \frac{s}{2}} \\ & + (P_1(c + \frac{s}{2}) - P_1(-c - \frac{s}{2})) \frac{s^2}{12} \end{aligned}$$

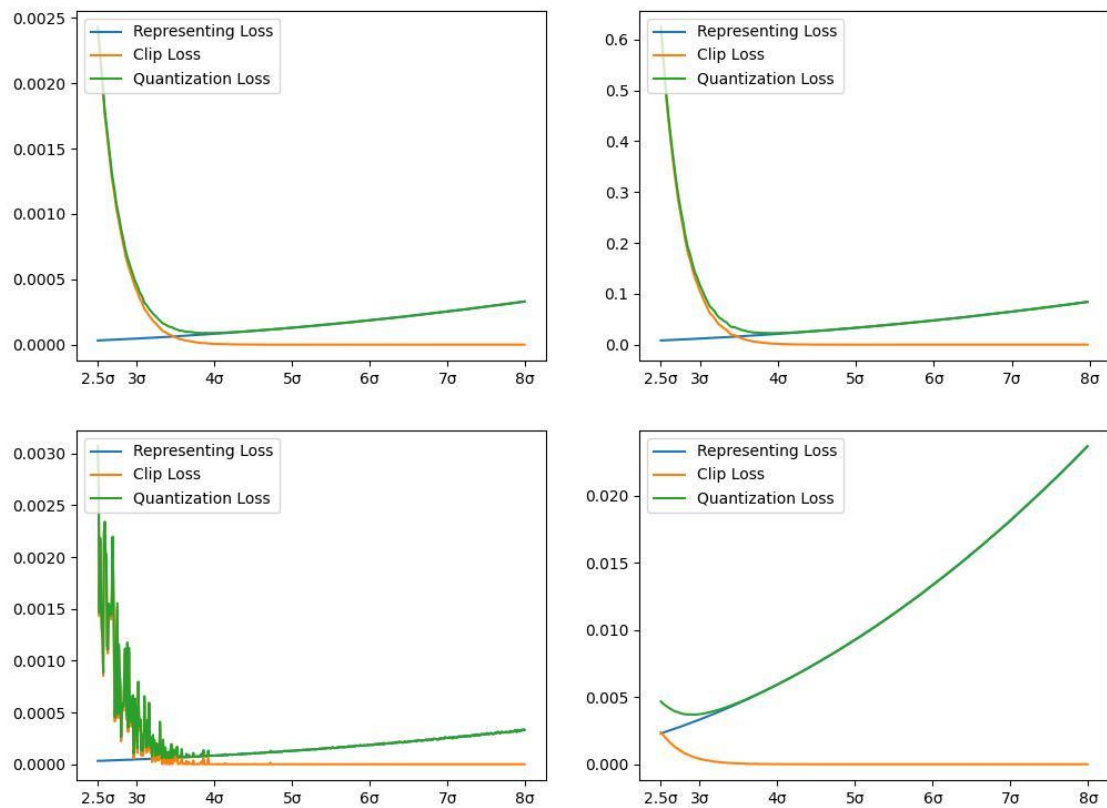
$$E\{L_{MSE}\} = \frac{s^3}{12} * \frac{2 \int_0^{k\sigma + \frac{s}{2}} \varphi(x) dx}{s} + \sqrt{\frac{\pi}{2}} (k^2 \sigma^2 + 1) (1 - erf(\frac{k\sigma + \frac{s}{2}}{\sqrt{2}})) + (k\sigma - \frac{s}{2}) e^{-\frac{(k\sigma + \frac{s}{2})^2}{2}}$$

量化误差分析 - Bernard Widrow公式

$$E\{L_{MSE}\} = (x - c)^2 P_1(x) - 2((x - c)P_2(x) - P_3(x)) \Big|_{c+\frac{s}{2}}^{inf} + (x + c)^2 P_1(x) - 2((x + c)P_2(x) - P_3(x)) \Big|_{-inf}^{-c-\frac{s}{2}} \\ + (P_1(c + \frac{s}{2}) - P_1(-c - \frac{s}{2})) \frac{s^2}{12}$$

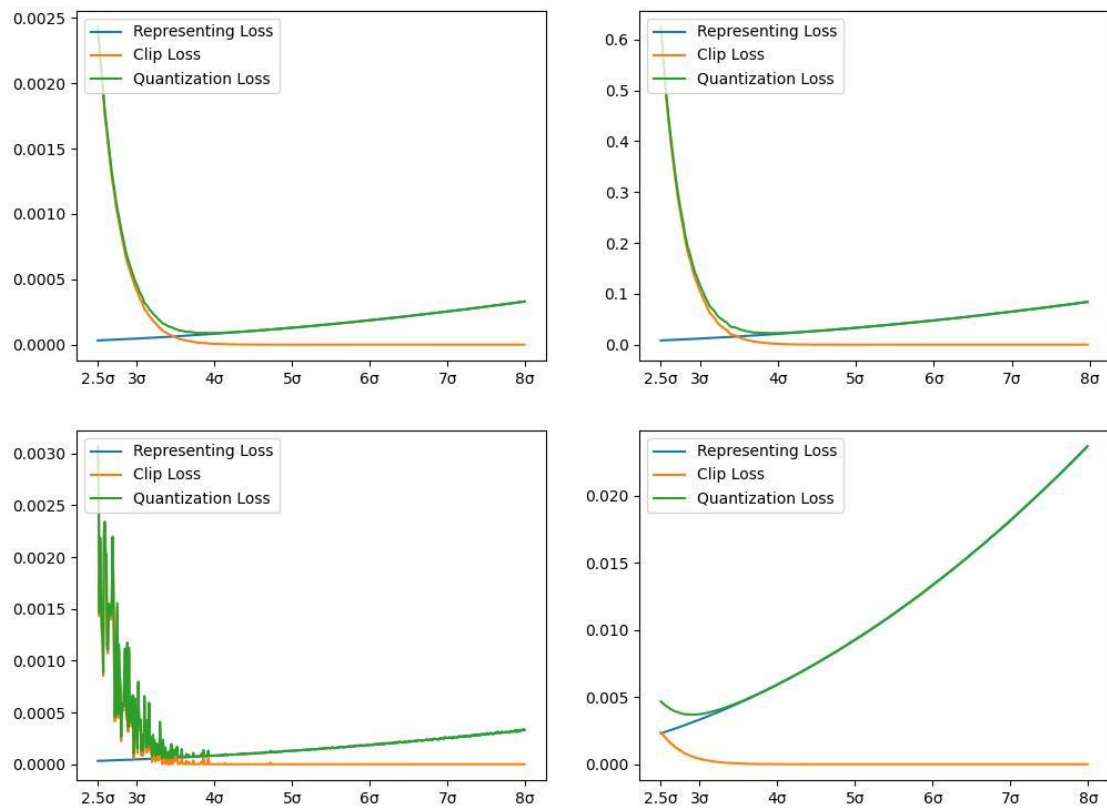
- 该等式是量化核心成果之一，该等式表明了量化误差是正负截断误差与表示误差的累计和。
- 对于高斯分布而言，截断误差随着截断值增长而指数级收敛；表示误差随着截断值增长而增长，增长速度为二次方级。
- 在截断值不变的情况下，表示误差随scale的增长而增长，增长速度为二次方级。

量化误差分析 - 截断误差C与表示误差R的关系



- 图左上：8bit定点 $\sigma=1$ 时的C-R关系图
- 图右上：8bit定点 $\sigma=16$ 时的C-R关系图
- 图左下：8bit定点 $\sigma=1$ 时的C-R关系图（只有4096个样本点）
- 图右上：4bit定点 $\sigma=1$ 时的C-R关系图

量化误差分析 - 截断误差C与表示误差R的关系



结论：

- 由于C的指数级收敛，在C-R关系图中R可以认为是一个线性增长的函数；量化参数选择实际上是C-R函数关系决定的。
- 对于不同的sigma而言，C-R关系保持稳定，可以选取相同的 $k \cdot \sigma$ 作为截断点。
- 样本量较少时估计的方差较大，估计不稳定，方差主要来自于大值的阶段误差。
- 4bit与8bit的C-R性质不同，说明了4bit优化与8bit优化存在差异。

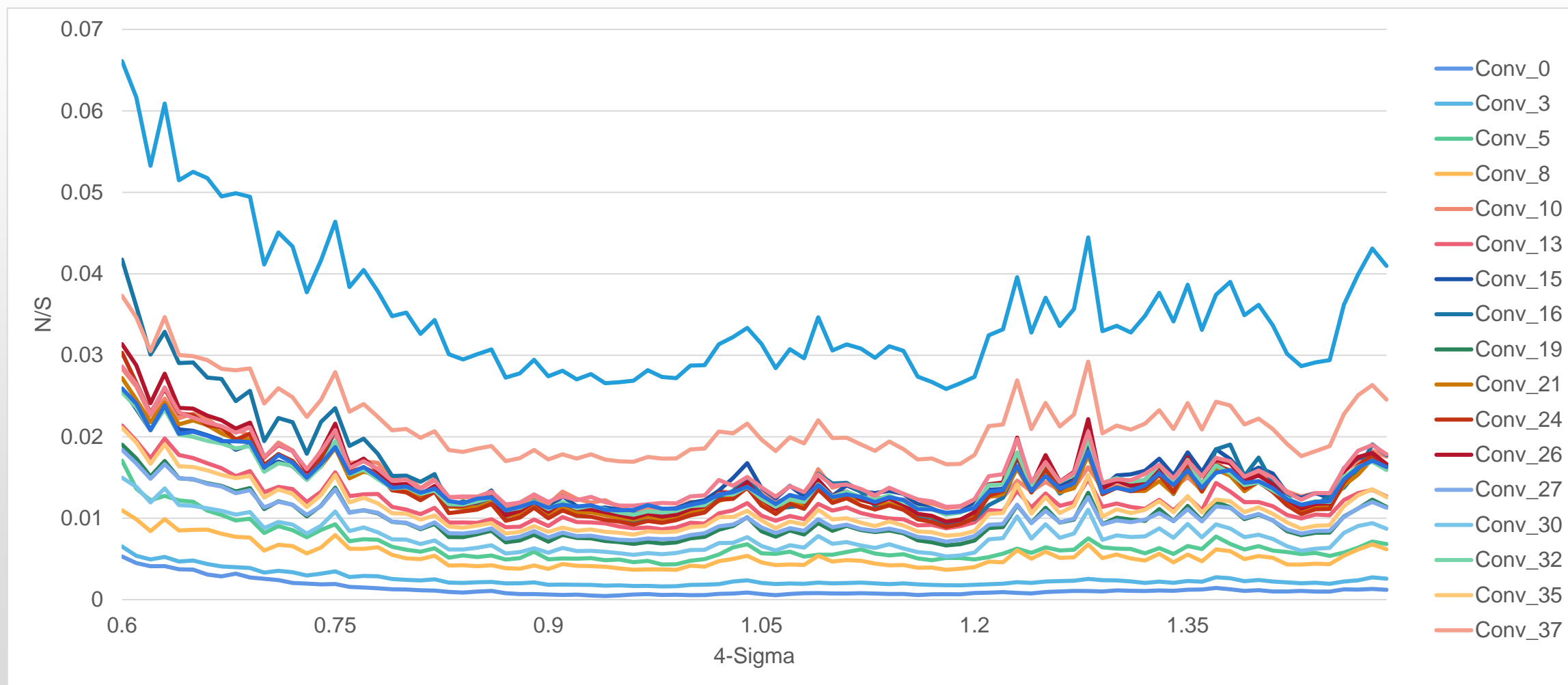
最优估计存在的问题

$$E\{L_{MSE}\} =$$

$$\begin{aligned} & (x-c)^2 P_1(x) - 2((x-c)P_2(x) - P_3(x)) \Big|_{c+\frac{s}{2}}^{inf} + (x+c)^2 P_1(x) - 2((x+c)P_2(x) - P_3(x)) \Big|_{-inf}^{-c-\frac{s}{2}} \\ & + (P_1(c+\frac{s}{2}) - P_1(-c-\frac{s}{2})) \frac{s^2}{12} \end{aligned}$$

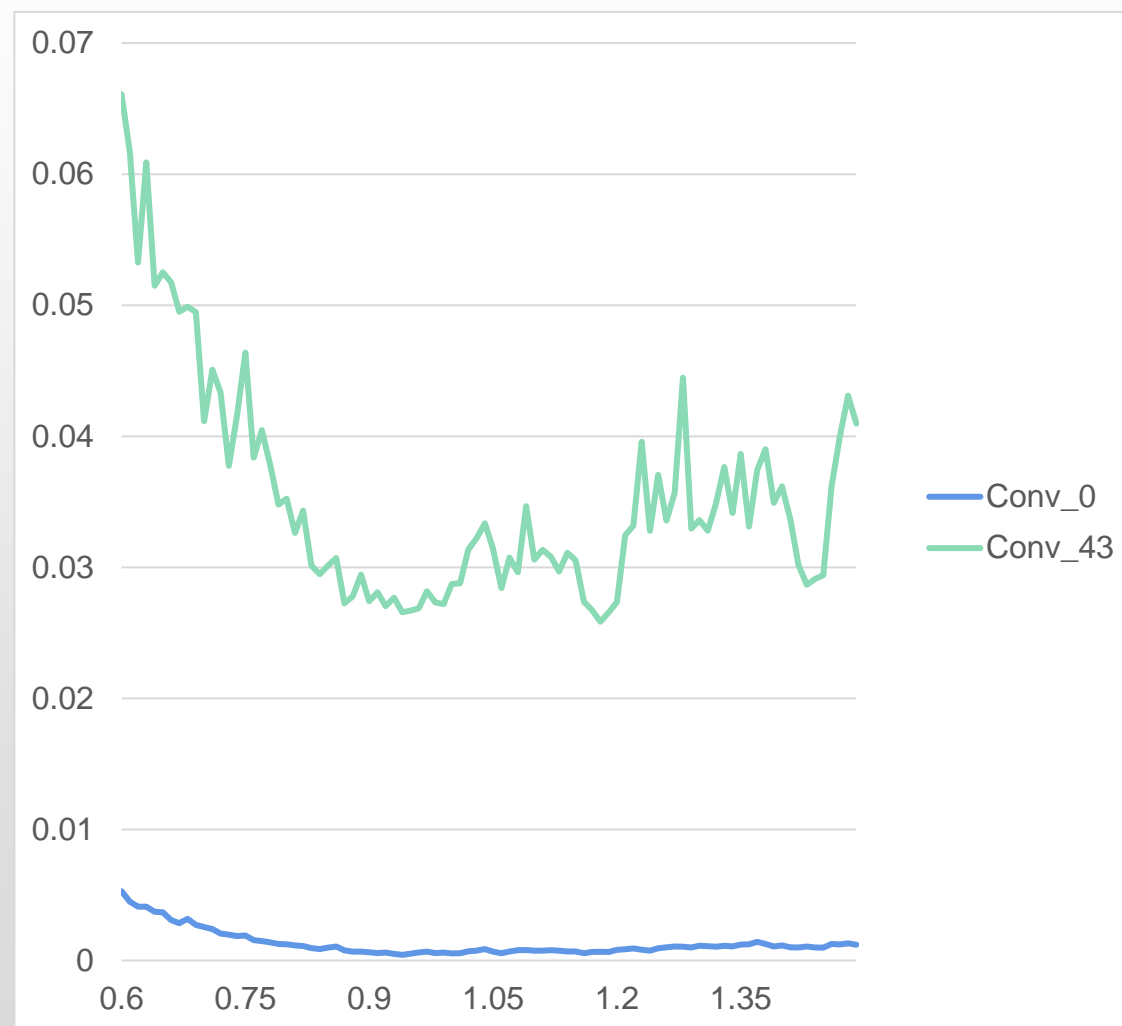
- 最优截断要求pdf的三阶积分，并求导令上式为0，对于大部分分布而言，无法顺利求得解析解。
- 同时在很多情况下，局部的MSE最优并不是全局MSE最优的。
- 数据量小时，估计的方差很大。

局部最优截断与全局最优截断



局部最优截断与全局最优截断

- 在resnet18的结果中我们不难发现：conv0的误差（局部量化误差）与conv43的误差（全局量化误差）之间的关系难以被建模，显然局部最优截断难以保证全局最优。



枚举最优截断

- 方案4：枚举截断值：

$$E\{L_{MSE}\} = \int_{-c-\frac{s}{2}}^{c+\frac{s}{2}} (x' - x)^2 \varphi(x) dx + \int_{c+\frac{s}{2}}^{inf} (x - k\sigma)^2 \varphi(x) dx + \int_{-inf}^{-c-\frac{s}{2}} (x + c)^2 \varphi(x) dx$$

- 考虑实际求解上述方程的最小值存在困难，可以直接改为直接用样本计算MSE，枚举截断值并从中取优。
 - 初始化枚举候选点，例如C=[0.1, 0.2, 0.3, 0.4, 0.5 ... 20.0]
 - 从C任取一个候选点 c' ，带入上式，求 $L_{MSE}|_{c'} = \sum (x'_i - x_i)^2$
 - 寻找一个最合适的 c^* ，使得 $L_{MSE}|_{c^*}$ 最小，作为截断值
- 算法复杂度为O(NM)，其中N为元素个数，M为候选点个数。
- 实践中我们使用直方图统计优化该算法，经过优化的算法复杂度为O(M)。

枚举最优截断

- 是否存在其他损失函数可以用来确定截断值？例如寻找最优SNR，最优KL散度？

$$L_{KL} = \int_{-inf}^{inf} p(x) \log \frac{p(x)}{q(x)}, \quad \text{where } q(x) = \frac{P(u_i) - P(l_i)}{s}$$

- 此处额外定义 $l_1 = -inf, u_{254} = inf$
 1. 初始化枚举候选点，例如 $C = [0.1, 0.2, 0.3, 0.4, 0.5 \dots 20.0]$
 2. 从C任取一个候选点 c' ，带入上式，求 $L_{KL}|_{c'} = \sum Count(l_i, u_i, X) \log \left(\frac{Count(l_i, u_i, X)}{Count(l_i, u_i, X')} \right)$
 3. 寻找一个最合适的 c^* ，使得 $L_{KL}|_{c'}$ 最小，作为截断值
- 算法复杂度为 $O(NM)$ ，其中N为元素个数，M为候选点个数。
- 实践中我们使用直方图统计优化该算法，经过优化的算法复杂度为 $O(M)$ 。
- $Count(l_i, u_i, X)$ ：统计X中有多少个元素落入区间 $(l_i, u_i]$

基于梯度优化的截断值

- 方案5：梯度优化截断值

- 我们定义量化函数 $Q(x, s) = \begin{cases} [\frac{x}{s}], & else \\ 127, & \frac{x}{s} > 127.5 \\ -127, & \frac{x}{s} < -127.5 \end{cases}$, 额外定义截断值 $c = s * 127.5$

- 我们定义反量化函数 $DQ(x, s) = x * s$

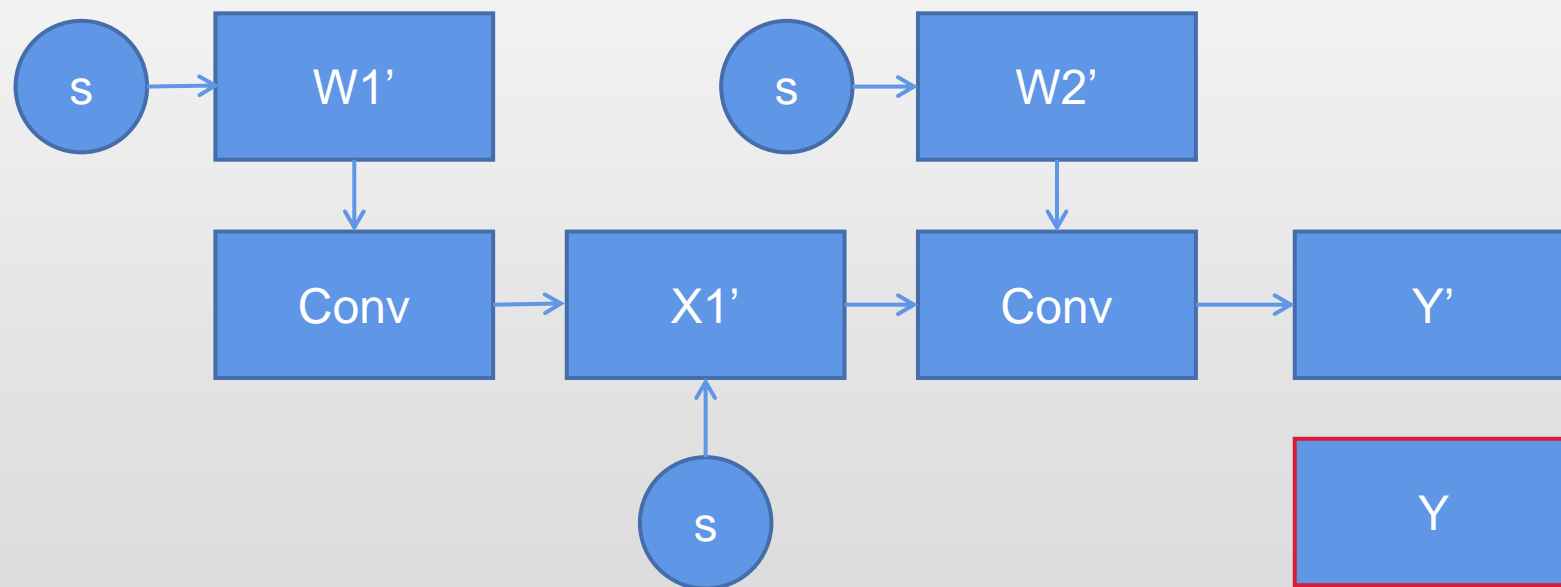
- 我们定义x的量化值为 $x'(s) \triangleq DQ(Q(x, s), s)$

- 思考：为了解决全局最优的问题，能否使用梯度优化 s ？

$$\frac{dx'}{ds} = ?$$

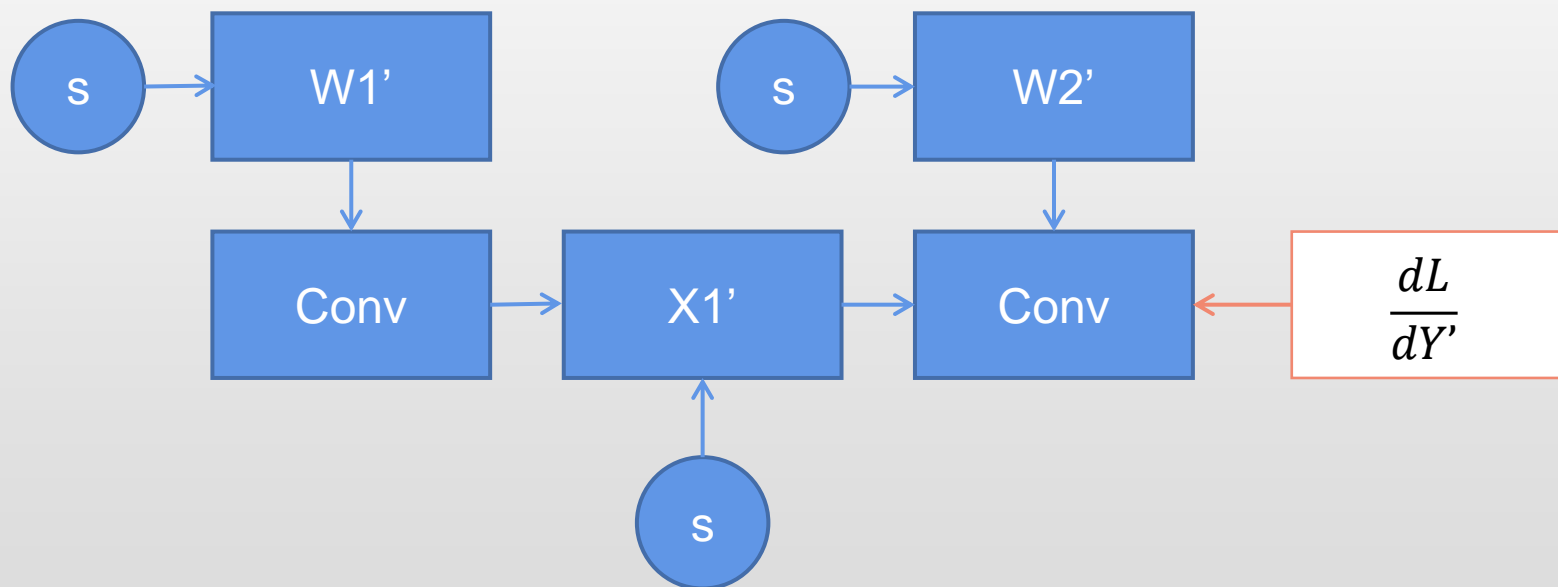
基于梯度优化的截断值

- 对于任意一个网络而言，我们总是能收集到量化后的网络输出 Y' ，以及量化前的网络输出 Y 。
- 因此我们可以直接计算 $L = \text{MSE}(Y', Y)$ ，形成损失函数并求导



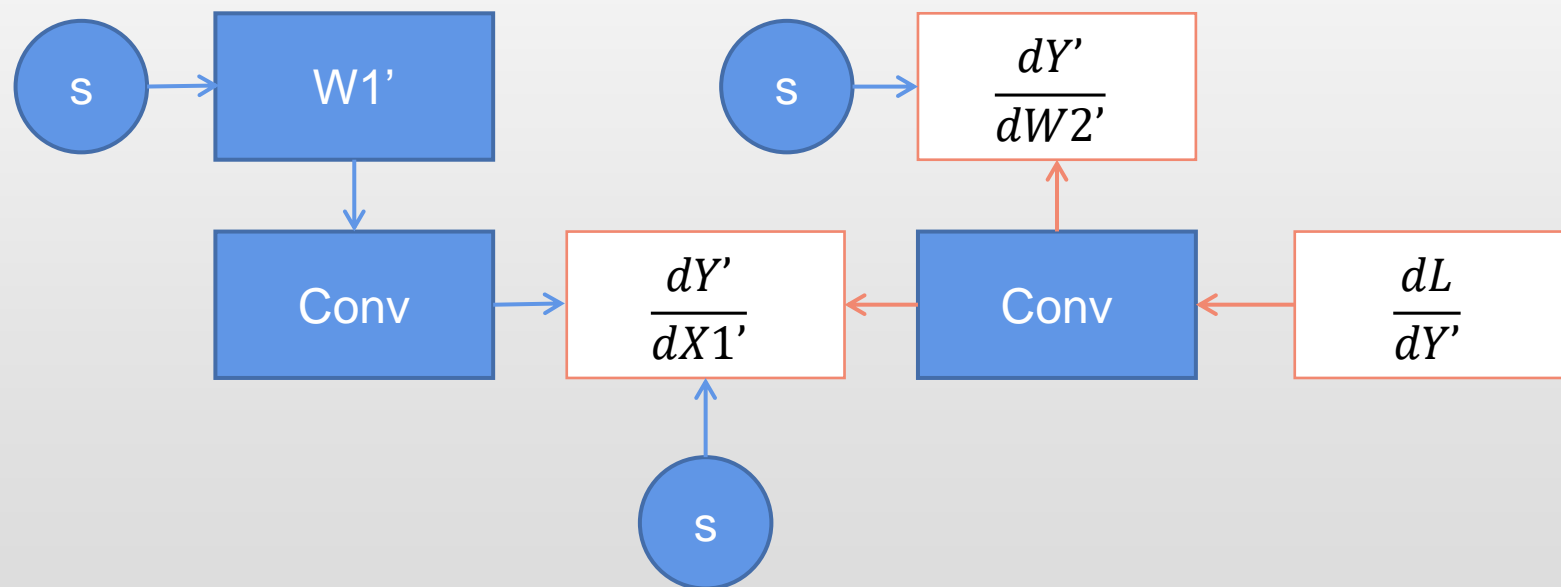
基于梯度优化的截断值

- 对于任意一个网络而言，我们总是能收集到量化后的网络输出 Y' ，以及量化前的网络输出 Y 。
- 因此我们可以直接计算 $L = \text{MSE}(Y', Y)$ ，形成损失函数并求导



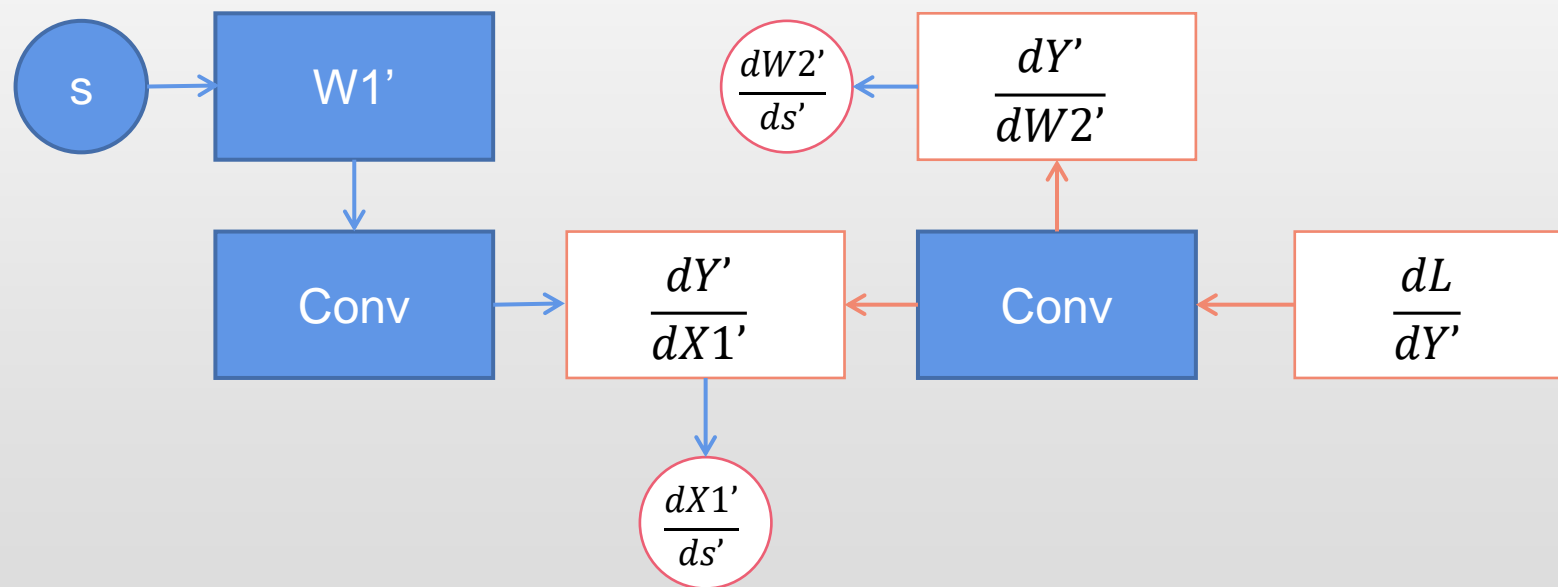
基于梯度优化的截断值

- 对于任意一个网络而言，我们总是能收集到量化后的网络输出 Y' ，以及量化前的网络输出 Y 。
- 因此我们可以直接计算 $L = \text{MSE}(Y', Y)$ ，形成损失函数并求导



基于梯度优化的截断值

- 问题：导数 $\frac{dW2'}{ds'}$ 怎么求？



基于梯度优化的截断值

- 问题：导数 $\frac{dW2'}{ds'}$ 怎么求？

- 我们定义量化函数 $Q(x, s) = \begin{cases} [\frac{x}{s}], & else \\ 127, & \frac{x}{s} > 127.5 \\ -127, & \frac{x}{s} < -127.5 \end{cases}$, 额外定义截断值 $c = s * 127.5$
- 我们定义反量化函数 $DQ(x, s) = x * s$
- 我们定义x的量化值为 $x'(s) \triangleq DQ(Q(x, s), s)$

基于梯度优化的截断值

- 问题：导数 $\frac{dW2'}{ds'}$ 怎么求？

- $$w'(s) = \begin{cases} w * [\frac{x}{s}], & else \\ w * 127, & \frac{w}{s} > 127.5 \\ -w * 127, & \frac{w}{s} < -127.5 \end{cases}$$

- $$dw'(s) = \begin{cases} ds * [\frac{w}{s}] + s * d[\frac{w}{s}], & else \\ ds * 127, & \frac{w}{s} > 127.5 \\ -ds * 127, & \frac{w}{s} < -127.5 \end{cases}$$

- $d[\frac{w}{s}] = 0$

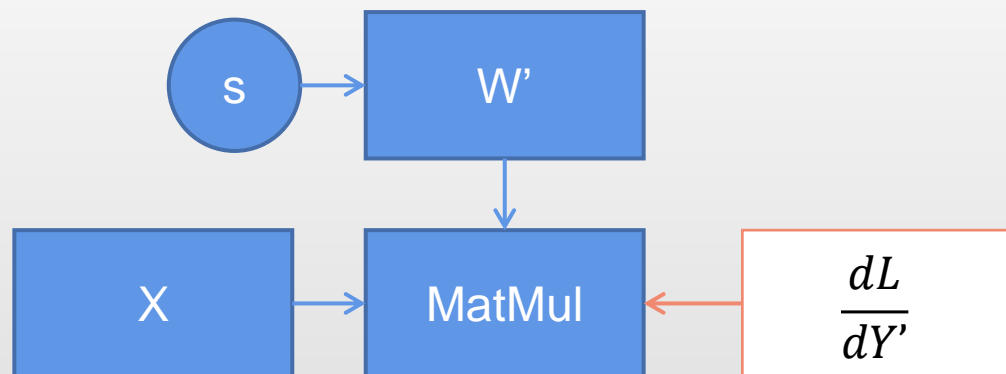
习题

- 已知 $Y = W'X$, 其中 W, X 均是三阶方阵

- $X = \begin{bmatrix} 1 & 4 & 5 \\ 1 & -2 & 3 \\ 0 & 3 & 0 \end{bmatrix}, W = \begin{bmatrix} 3 & 3 & 5 \\ 0 & 4 & -3 \\ 1 & 1 & 1 \end{bmatrix}$

- $w'(s) = \begin{cases} w * [\frac{x}{s}], & else \\ w * 127, & \frac{w}{s} > 127.5 \\ -w * 127, & \frac{w}{s} < -127.5 \end{cases}$

- 求 s 关于 $(\text{sum}(Y) - 10)^2$ 在 $s = 0.02$ 处的导数



相关代码

- 访问 <https://github.com/openppl-public/ppq/tree/master/ppq>