

Proof of LQER with an Assumption to be Verified

| | |
|---------------------------------------------------------|---|
| Proof of LQER with an Assumption to be Verified | 1 |
| Background | 1 |
| Geroge's Proof with Assumption | 1 |
| An Example of Statistic Profile | 2 |
| What if the Assumption does not hold? | 3 |
| Cheng's Try without the Assumption | 4 |
| Numerical Stability of Matrix Square Root and SVD | 5 |

Background

- Input vector of linear layer: $\mathbf{x} = [x_1, x_2, \dots, x_n]$, where n is the hidden size.
- linear layer: $\mathbf{y} = \mathbf{x}W$
- Approximated linear layer: $\tilde{\mathbf{y}} = \mathbf{x}(\tilde{W} + C)$, where C is a low-rank approximation of quantization error, and $\tilde{W} = \text{quantize}(W)$

Geroge's Proof with Assumption

We target minimizing the expectation of L2 norm for output vectors, *i.e.*, $\min \mathbb{E}_{\mathbf{y}} \{ \|\tilde{\mathbf{y}} - \mathbf{y}\|_2^2 \}$.

We have

$$\mathbb{E}_{\mathbf{y}} \{ \|\tilde{\mathbf{y}} - \mathbf{y}\|_2^2 \} = \mathbb{E}_{\mathbf{x}} \{ \|\mathbf{x}(\tilde{W} - W + C)\|_2^2 \} \quad (1)$$

where $\mathbb{E}\{\cdot\}$ stands for expectation.

Let $M = \tilde{W} - W + C = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \vdots \\ \mathbf{m}_n \end{bmatrix}$ and put it into RHS of (1).

$$\begin{aligned} \mathbb{E}_{\mathbf{y}} \{ \|\tilde{\mathbf{y}} - \mathbf{y}\| \} &= \mathbb{E}_{\mathbf{x}} \{ \|\mathbf{x}M\|_2^2 \} \\ &= \mathbb{E}_{\mathbf{x}} \left\{ \left\| [x_1, x_2, \dots, x_n] \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \vdots \\ \mathbf{m}_n \end{bmatrix} \right\|_2^2 \right\} \\ &= \mathbb{E}_{\mathbf{x}} \left\{ \left\| \sum_{i=1}^n x_i \mathbf{m}_i \right\|_2^2 \right\} \\ &= \mathbb{E}_{\mathbf{x}} \left\{ \left(\sum_{i=1}^n x_i \mathbf{m}_i \right) \left(\sum_{j=1}^n x_j \mathbf{m}_j^T \right) \right\} \\ &= \mathbb{E}_{\mathbf{x}} \left\{ \sum_{i=1}^n \sum_{j=1}^n x_i x_j \mathbf{m}_i \mathbf{m}_j^T \right\} \quad \text{⚓} \\ &= \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}_{\mathbf{x}} \{ x_i x_j \mathbf{m}_i \mathbf{m}_j^T \} \\ &= \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}_{\mathbf{x}} \{ x_i x_j \} \mathbf{m}_i \mathbf{m}_j^T \end{aligned} \quad (2)$$

Assumption:

$$\mathbb{E}\{x_i x_j\} = 0 \text{ for } i \neq j \quad (3)$$

Then the RHS of (2) becomes

$$\mathbb{E}_{\mathbf{y}}\{\|\tilde{\mathbf{y}} - \mathbf{y}\|\} = \sum_{i=1}^n \mathbb{E}\{x_i^2\} \mathbf{m}_i \mathbf{m}_i^T \quad (4)$$

💡 One interpretation of (3) is that the i -th activation dim x_i is zero-mean, and independent from the j -th dim ($i \neq j$).

If we assign diagonal matrix $S = \text{diag}\left(\sqrt{\mathbb{E}\{x_1^2\}}, \sqrt{\mathbb{E}\{x_2^2\}}, \dots, \sqrt{\mathbb{E}\{x_n^2\}}\right)$, the LHS of (4) becomes

$$\begin{aligned} \mathbb{E}_{\mathbf{y}}\{\|\tilde{\mathbf{y}} - \mathbf{y}\|\} &= \text{Trace}(SMM^T S^T) \\ &= \|SM\|_F^2 \end{aligned} \quad (5)$$

where $\|\cdot\|_F^2$ denotes Frobenius norm.

Given (5), our target now is

$$\begin{aligned} \min \mathbb{E}_{\mathbf{y}}\{\|\tilde{\mathbf{y}} - \mathbf{y}\|_2^2\} &= \min \|SM\|_F^2 \\ &= \min \|S(\tilde{W} - W + C)\|_F^2 \end{aligned} \quad (6)$$

If we assign $A = S(W - \tilde{W})$ and $\tilde{A} = SC$, the RHS of (6) becomes

$$\min \mathbb{E}_{\mathbf{y}}\{\|\tilde{\mathbf{y}} - \mathbf{y}\|_2^2\} = \min \|\tilde{A} - A\|_F^2 \quad (7)$$

According to [Eckart-Young theorem](#), the best rank k approximation to A (noted as \tilde{A}_k) is

$$\tilde{A}_k = U_k \Sigma_k V_k^T \quad (8)$$

where U_k, Σ_k, V_k is the rank- k SVD of A , i.e., $A = S(W - \tilde{W}) = U\Sigma V^T$.

Therefore, the closed-form solution of C based on assumption (3) is

$$C = S^{-1}\tilde{A}_k = S^{-1}U_k \Sigma_k V_k^T \quad (9)$$

For implementation, we assign two low-rank matrices $S^{-1}U_k$ and $\Sigma_k V_k^T$ to save FLOPs.

An Example of Statistic Profile

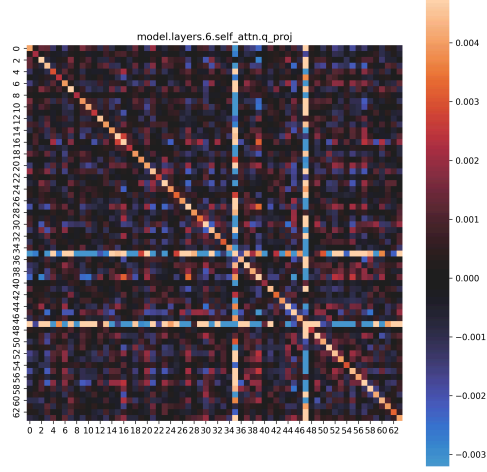
We can calculate the [auto-correlation](#) matrix for $\mathbf{x} = [x_1, x_2, \dots, x_n]$:

$$R_{\mathbf{xx}} = \begin{pmatrix} \mathbb{E}(x_1 x_1) & \mathbb{E}(x_1 x_2) & \dots & \mathbb{E}(x_1 x_n) \\ \mathbb{E}(x_2 x_1) & \mathbb{E}(x_2 x_2) & \dots & \mathbb{E}(x_2 x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{E}(x_n x_1) & \mathbb{E}(x_n x_2) & \dots & \mathbb{E}(x_n x_n) \end{pmatrix} \quad (10)$$

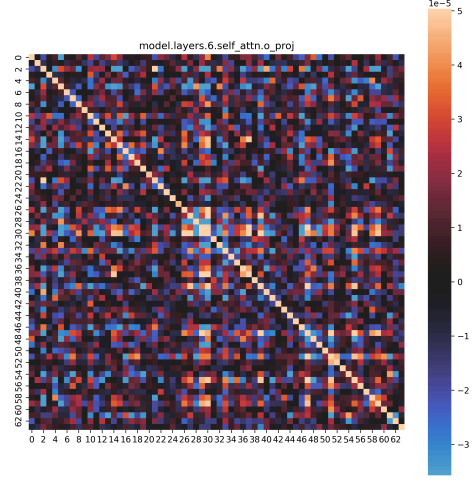
where $\mathbb{E}(x_i x_j)$ is the correlation between \mathbf{x} 's i -th dim and j -th dim.

Figure 1 is the profiled first 64 dim of R_{xx} for the linear layers of 6-th decoder in [TinyLlama-1.1B](#). Color black mean 0 value.

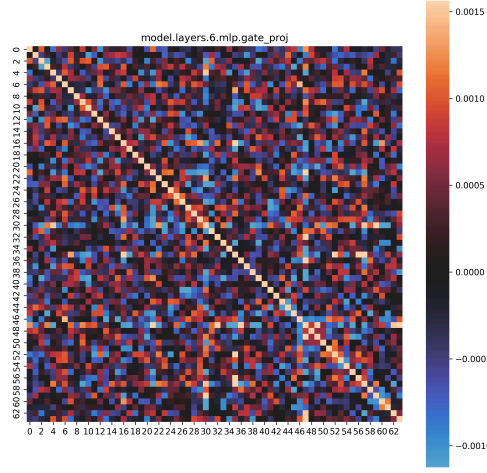
It is observed that our assumption is reasonable but looks a bit strong (?).



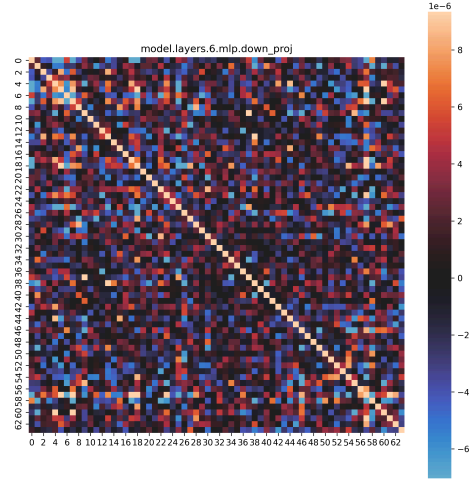
(a) The input of q/k/v_proj layer of 6-th decoder



(b) The input of o_proj layer of 6-th decoder



(c) The input of gate/up_proj layer of 6-th decoder



(d) The input of down_proj layer of 6-th decoder

What if the Assumption does not hold?

George hypothesizes that if the assumption (3) does not hold, the closed-form of S might be $S = R_{\mathbf{x}\mathbf{x}}^{\frac{1}{2}}$, where $R_{\mathbf{x}\mathbf{x}} = \mathbb{E}\{\mathbf{x}\mathbf{x}^T\}$. Then (2) may become $\| R_{\mathbf{x}\mathbf{x}}^{\frac{1}{2}} \odot M \|_F^2$ where \odot is elementwise-multiply.

Cheng's Try without the Assumption

Cheng: I tried the following to push the proof.

If we continue at \S in (2), we have:

$$\begin{aligned}
 E_{\mathbf{y}}\{\|\tilde{\mathbf{y}} - \mathbf{y}\|\} &= E_{\mathbf{x}}\left\{\sum_{i=1}^n \sum_{j=1}^n x_i x_j \mathbf{m}_i \mathbf{m}_j^T\right\} \\
 &= E_{\mathbf{x}}\left\{\mathbf{e} \begin{pmatrix} x_1 x_1 & x_1 x_2 & \dots & x_1 x_n \\ x_2 x_1 & x_2 x_2 & \dots & x_2 x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_n x_1 & x_n x_2 & \dots & x_n x_n \end{pmatrix} \odot \begin{pmatrix} \mathbf{m}_1 \mathbf{m}_1^T & \mathbf{m}_1 \mathbf{m}_2^T & \dots & \mathbf{m}_1 \mathbf{m}_n^T \\ \mathbf{m}_2 \mathbf{m}_1^T & \mathbf{m}_2 \mathbf{m}_2^T & \dots & \mathbf{m}_2 \mathbf{m}_n^T \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{m}_n \mathbf{m}_1^T & \mathbf{m}_n \mathbf{m}_2^T & \dots & \mathbf{m}_n \mathbf{m}_n^T \end{pmatrix} \mathbf{e}^T\right\} \quad \S \quad (11) \\
 &= E_{\mathbf{x}}\{\mathbf{e} \cdot ((\mathbf{x}^T \mathbf{x}) \odot (MM^T)) \cdot \mathbf{e}^T\}
 \end{aligned}$$

where $\mathbf{e} = \underbrace{[1, 1, \dots, 1]}_n$ is a row vector of n ones, and \odot is elementwise multiply. \S in (11) is based on this [StackExchange page](#).

Using the following [property of Hadamard product](#)

- For vectors \mathbf{x} and \mathbf{y} , and corresponding diagonal matrices $D_{\mathbf{x}}$ and $D_{\mathbf{y}}$ with these vectors as their main diagonals, the following identity holds:^{[1]:479}

$$\mathbf{x}^* (A \odot B) \mathbf{y} = \text{tr}(D_{\mathbf{x}}^* A D_{\mathbf{y}} B^T),$$

where \mathbf{x}^* denotes the [conjugate transpose](#) of \mathbf{x} . In particular, using vectors of ones, this shows that the sum of all elements in the Hadamard product is the [trace](#) of AB^T where superscript T denotes the [matrix transpose](#), that is, $\text{tr}(AB^T) = \mathbf{1}^T (A \odot B) \mathbf{1}$. A related result for square A and B , is that the row-sums of their Hadamard product are the diagonal elements of AB^T :^[8]

$$\sum_i (A \odot B)_{ij} = (B^T A)_{jj} = (AB^T)_{ii}.$$

Similarly,

$$(\mathbf{y} \mathbf{x}^*) \odot A = D_{\mathbf{y}} A D_{\mathbf{x}}^*$$

Furthermore, a Hadamard matrix-vector product can be expressed as:

$$(A \odot B) \mathbf{y} = \text{diag}(A D_{\mathbf{y}} B^T)$$

where $\text{diag}(M)$ is the vector formed from the diagonals of matrix M .

The LHS of (11) becomes

$$\begin{aligned}
 E_{\mathbf{y}}\{\|\tilde{\mathbf{y}} - \mathbf{y}\|\} &= E_{\mathbf{x}}\{\text{Trace}(\mathbf{x}^T \mathbf{x} M M^T)\} \\
 &= \text{Trace}(E_{\mathbf{x}}\{\mathbf{x}^T \mathbf{x} M M^T\}) \quad \S \\
 &= \text{Trace}(E_{\mathbf{x}}\{\mathbf{x}^T \mathbf{x}\} M M^T) \\
 &= \text{Trace}(R_{\mathbf{xx}} M M^T) \\
 &= \text{Trace}\left(R_{\mathbf{xx}}^{\frac{1}{2}} M M^T R_{\mathbf{xx}}^{\frac{1}{2}}\right) \quad \S \\
 &= \text{Trace}\left(R_{\mathbf{xx}}^{\frac{1}{2}} M M^T \left(R_{\mathbf{xx}}^{\frac{1}{2}}\right)^T\right) \\
 &= \|R_{\mathbf{xx}}^{\frac{1}{2}} M\|_F^2
 \end{aligned} \quad (12)$$

\S in (12) is based on [this page](#).

\S in (12) is based on the following clues:

-
- $R_{\mathbf{xx}}$ is symmetric and positive-semidefinite. Refer to [wiki page](#).
 - $M M^T$ is positive-semidefinite. Refer to [StackExchange page](#)
 - For two positive-semidefinite matrices A and B , we have $\text{Trace}(AB) = \text{Trace}\left(A^{\frac{1}{2}} B A^{\frac{1}{2}}\right)$: refer to [wiki page](#)

Trace [\[edit\]](#)

The diagonal entries m_{ii} of a positive-semidefinite matrix are real and non-negative. As a consequence the [trace](#), $\text{tr}(M) \geq 0$. Furthermore,^[13] since every principal sub-matrix (in particular, 2-by-2) is positive semidefinite,

$$|m_{ij}| \leq \sqrt{m_{ii}m_{jj}} \quad \forall i, j$$

and thus, when $n \geq 1$,

$$\max_{i,j} |m_{ij}| \leq \max_i m_{ii}$$

An $n \times n$ Hermitian matrix M is positive definite if it satisfies the following trace inequalities:^[14]

$$\text{tr}(M) > 0 \quad \text{and} \quad \frac{(\text{tr}(M))^2}{\text{tr}(M^2)} > n - 1.$$

Another important result is that for any M and N positive-semidefinite matrices, $\text{tr}(MN) \geq 0$. This follows by writing $\text{tr}(MN) = \text{tr}(M^{\frac{1}{2}} N M^{\frac{1}{2}})$.

The matrix $M^{\frac{1}{2}} N M^{\frac{1}{2}}$ is positive-semidefinite and thus has non-negative eigenvalues, whose sum, the trace, is therefore also non-negative.

- R_{xx} has a symmetric and positive semidefinite square root: $R_{\text{xx}} = R_{\text{xx}}^{\frac{1}{2}} R_{\text{xx}}^{\frac{1}{2}}$. Refer to [wiki page](#)

Positive semidefinite matrices [\[edit\]](#)

See also: [Positive definite matrix § Decomposition](#)

A symmetric real $n \times n$ matrix is called [positive semidefinite](#) if $x^T A x \geq 0$ for all $x \in \mathbb{R}^n$ (here x^T denotes the [transpose](#), changing a column vector x into a row vector). A square real matrix is positive semidefinite if and only if $A = B^T B$ for some matrix B . There can be many different such matrices B . A positive semidefinite matrix A can also have many matrices B such that $A = B B$. However, A always has precisely one square root B that is [positive semidefinite and symmetric](#). In particular, since B is required to be symmetric, $B = B^T$, so the two conditions $A = B B$ or $A = B^T B$ are equivalent.

Now our target becomes

$$\begin{aligned} \min_{\mathbf{y}} E_{\mathbf{y}} \{ \|\tilde{\mathbf{y}} - \mathbf{y}\|_2^2 \} &= \min_{\mathbf{y}} \| R_{\text{xx}}^{\frac{1}{2}} M \|_F^2 \\ &= \min_{\mathbf{y}} \| R_{\text{xx}}^{\frac{1}{2}} (\tilde{W} - W + C) \|_F^2 \end{aligned} \quad (13)$$

If we assign $A = R_{\text{xx}}^{\frac{1}{2}} (W - \tilde{W})$ and $\tilde{A} = R_{\text{xx}}^{\frac{1}{2}} C$, the RHS of (13) becomes

$$\min_{\mathbf{y}} E_{\mathbf{y}} \{ \|\tilde{\mathbf{y}} - \mathbf{y}\|_2^2 \} = \min_{\mathbf{y}} \| \tilde{A} - A \|_F^2 \quad (14)$$

According to [Eckart–Young theorem](#), the best rank k approximation to A (noted as \tilde{A}_k) is

$$\tilde{A}_k = U_k \Sigma_k V_k^T \quad (15)$$

Therefore, the closed-form solution of C without the assumption is

$$C = \left(R_{\text{xx}}^{\frac{1}{2}} \right)^{-1} \tilde{A}_k = \left(R_{\text{xx}}^{\frac{1}{2}} \right)^{-1} U_k \Sigma_k V_k^T \quad (16)$$

For implementation, we assign two low-rank matrices $\left(\left(R_{\text{xx}}^{\frac{1}{2}} \right)^{-1} U_k \right)$ and $(\Sigma_k V_k^T)$ to save FLOPs.

? A concern here is if $R_{\text{xx}}^{\frac{1}{2}}$ is always invertible? We assume in practice it is invertible for a trained network, but we need a backup plan if it is not (Probably add perturbation to $R_{\text{xx}}^{\frac{1}{2}}$, or fall back to the diagonal S case)

Numerical Stability of Matrix Square Root and SVD

- The numerical stability of this algorithm is not ensured and has not been investigated.
- Current empirical implementation uses FP32 to compute R_{xx} , use FP64 to accumulated R_{xx} , and use FP32 to solve the matrix square root and SVD. The evaluation of the quantized network may use FP32, or BF16 to align with MXINT's 8-bit exponent.