

---

# Typst template for mathematical papers

---

Léo Burgund

April 25, 2025

## 1 Nomenclature

### 1.1 Dimensions

- $b$  Mini-batch size
- $d_e$  Embedding dimension
- $d_s$  Sequence length
- $d_k$  Query/Keys dimension
- $d_v$  Value dimension
- $h$  Number of heads

### 1.2 Matrix

We will first place ourselves in the case where  $b = 1$ .

In the case of multi head attention, for each head  $i = 1, \dots, h$ , we have:

- Input  $X \in \mathbb{R}^{d_s \times d_e}$
- $W_{Q_i} \in \mathbb{R}^{d_e \times \frac{d_k}{h}}$ ,  $Q_i := XW_{Q_i} \in \mathbb{R}^{d_s \times \frac{d_k}{h}}$
- $W_{K_i} \in \mathbb{R}^{d_e \times \frac{d_k}{h}}$ ,  $K_i := XW_{K_i} \in \mathbb{R}^{d_s \times \frac{d_k}{h}}$
- $S_i := \frac{Q_i K_i^\top}{\sqrt{\frac{d_k}{h}}} \in \mathbb{R}^{d_s \times d_s}$
- $A_i := \text{softmax}_{\text{row}}(S)$
- $W_{V_i} \in \mathbb{R}^{d_e \times \frac{d_v}{h}}$ ,  $V_i := XW_{V_i} \in \mathbb{R}^{d_s \times \frac{d_v}{h}}$
- $H_i := A_i V_i \in \mathbb{R}^{d_s \times \frac{d_v}{h}}$ ,  $H = [H_1, \dots, H_h] \in \mathbb{R}^{d_s \times d_v}$
- $W_O \in \mathbb{R}^{d_v \times d_e}$
- Output  $Y := HW_O + X \in \mathbb{R}^{d_s \times d_e}$

**Remark 1.1.** The number of parameters to learn

$$\left( \underbrace{2 \left( d_e \frac{d_k}{h} \right)}_{W_{Q_i}, W_{K_i}} + \underbrace{d_e \frac{d_v}{h}}_{W_{V_i}} \right) h + \underbrace{d_v d_e}_{W_O}$$

is the same for any  $h \in \mathbb{N}_+^*$ .

**Remark 1.2.** We can easily consider the bias by augmenting the matrices:

$$X' = [X \mid \mathbf{1}] \in \mathbb{R}^{d_s \times (d_e + 1)}$$

$$H' = [H \mid \mathbf{1}] \in \mathbb{R}^{d_s \times (d_v + 1)}$$

And adding a row of parameters to  $W_{Q_i}, W_{K_i}, W_{V_i}, W_O$ . For example:

$$W'_{Q_i} = \begin{pmatrix} W_{Q_i} \\ (b^Q)^\top \end{pmatrix} \in \mathbb{R}^{(d_e + 1) \times \frac{d_k}{h}}.$$

## 2 Problem

We study the case where  $h = 1$ .

We are interested in growing the  $d_k$  dimension. We consider the first order approximation, using the functional gradient,

$$\mathcal{L}(f + \partial f(d\theta, d\mathcal{A})) = \mathcal{L}(f) + \langle \nabla_f \mathcal{L}(f), \partial f(\partial\theta, \partial\mathcal{A}) \rangle + o(\|\partial f(\partial\theta, \partial\mathcal{A})\|).$$

To avoid the softmax's non linearity, we will consider the gradient with respect to the matrix  $S$ , just before the softmax.

We then have

$$\mathcal{L}(S + \partial S) = \mathcal{L}(S) + \langle \nabla_S \mathcal{L}(S), \partial S \rangle + o(\|\partial S\|)$$

with

$$\partial S = X(W_Q + \partial W_Q)(W_K + \partial W_K)^\top X^\top - XW_Q W_K^\top X^\top.$$

We have the following optimization problem:

$$\arg \min_{\partial S} \langle \nabla_S \mathcal{L}(S), \partial S \rangle, \text{ such that } \|\partial S\| \leq \gamma$$

$$\arg \min_{\partial W_Q, \partial W_K} \left\| B - X(W_Q + \partial W_Q)(W_K + \partial W_K)^\top X^\top \right\|_F^2$$

$$\text{with } B := \nabla_S \mathcal{L}(S) + XW_Q W_K^\top X^\top$$

Which is a low rank regression limited by  $d_k$  (if  $d_k < d_e$ ).  $B$  is known.

We can approximate  $\underbrace{X(W_Q + \partial W_Q)}_{d_e \times d_k} \underbrace{(W_K + \partial W_K)^\top X^\top}_{d_k \times d_e}$  with a truncated SVD, taking the first  $d_k$  singular values.

If we want to grow the inner dimension of the attention matrix by  $p$  neurons, we can instead approximate by taking the first  $d_{k'} := d_k + p$  singular values.

Hence, instead of approximating a matrix  $\underbrace{(W_Q + \partial W_Q)}_{d_e \times d_k} \underbrace{(W_K + \partial W_K)^\top}_{d_k \times d_e}$ , we approximate

$$\underbrace{Z}_{d_e \times d_e} = \underbrace{\tilde{W}_Q}_{d_e \times (d_{k'})(d_{k'}) \times d_e} \underbrace{\tilde{W}_K^\top}_{(d_{k'}) \times d_e} = \begin{bmatrix} W_Q + \partial W_Q & \underbrace{\tilde{W}_Q}_{d_e \times p} \end{bmatrix} \begin{bmatrix} W_K + \partial W_K & \underbrace{\tilde{W}_K}_{d_e \times p} \end{bmatrix}^\top$$

with  $\text{rank}(Z) \leq d_{k'}$  (we make the hypothesis that  $d_{k'} < d_e$ ).

We then have the optimization problem

$$\arg \min_Z \|B - XZX^\top\|_F^2 \text{ subject to } \text{rank}(Z) \leq d_{k'}.$$

Which is a low rank regression problem, limited by  $d_{k'}$ .

Let  $f$  such that

$$f(Z) = \|B - XZX^\top\|_F^2,$$

$f$  is convex.

We have

$$\nabla_Z f = -2X^\top(B - XZX^\top)X,$$

so

$$\nabla_Z f = 0 \iff X^\top XZX^\top X = X^\top BX.$$

In the case where  $d_e \leq d_s$  and  $\text{rank}(X) = d_e$ , then  $X^\top X$  is non-singular, and we have the solution

$$Z^* = (X^\top X)^{-1} X^\top BX (X^\top X)^{-1}.$$

In the general case,

$$Z^* = X^+ B (X^+)^{\top},$$

with  $X^+ = (X^\top X)^{-1} X^\top$  the Moore-Penrose inverse.

If we had  $d_{k'} \geq d_e$ , we could use the trivial factorization  $\mathring{W}_Q = Z^*, \mathring{W}_K = I_{d_e}$ .

As most of the time  $d_{k'} < d_e$ , we have to approximate the factorization.

According to the Eckart–Young–Mirsky theorem, the best approximation  $Z_{k'}^*$  of  $X^+ B (X^+)^{\top}$  with  $\text{rank}(Z_{k'}^*) = d_{k'}$  is obtained with a truncated SVD.

Indeed, we have

$$Z^* = U \Sigma V^\top, \quad \Sigma = \text{diag}(\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{d_e}).$$

We keep the  $d_{k'}$  largest singular values

$$U_{k'} = [u_1, \dots, u_{d_{k'}}], \quad V_{k'} = [v_1, \dots, v_{d_{k'}}], \quad \Sigma_{k'} = \text{diag}(\sigma_1, \dots, \sigma_{d_{k'}}).$$

We get

$$Z_{k'}^* = U_{k'} \Sigma_{k'} V_{k'}^\top, \quad \text{rank}(Z_{k'}^*) = d_{k'}$$

$$\mathring{W}_Q^* = U_{k'} \Sigma_{k'}^{\frac{1}{2}}, \quad \mathring{W}_K^* = V_{k'} \Sigma_{k'}^{\frac{1}{2}}.$$

**Remark 2.1.**

$$\min_{\mathring{W}_Q, \mathring{W}_K} \|B - X \mathring{W}_Q \mathring{W}_K^\top X^\top\|_F^2 = \sum_{i > d_{k'}} \sigma_i^2, \quad \text{subject to } \text{rank}(\mathring{W}_Q \mathring{W}_K^\top) \leq d_{k'}$$

**Remark 2.2.** For implementation:

Keep the matrices apart, for example for the weight matrix of  $Q$  :

$$\dot{W}_Q = W'_Q + \partial W'_Q + W_Q^{\text{new}}$$

with (remind that  $d_{k'} = d_k + p$  )

$$W'_Q = \begin{bmatrix} w_1 & \dots & w_k & | & \mathbf{0}_1 & \dots & \mathbf{0}_p \end{bmatrix}_{d_e \times (d_k + p)}$$

$$\partial W'_Q = \begin{bmatrix} \partial w_1 & \dots & \partial w_k & | & \mathbf{0}_1 & \dots & \mathbf{0}_p \end{bmatrix}_{d_e \times (d_k + p)}$$

$$W_Q^{\text{new}} = \begin{bmatrix} \mathbf{0}_1 & \dots & \mathbf{0}_k & | & w_1^{\text{new}} & \dots & w_p^{\text{new}} \end{bmatrix}_{d_e \times (d_k + p)}$$

with any vector  $w \in \mathbb{R}^{d_e}$ , and  $\mathbf{0} \in \mathbb{R}^{d_e}$  the 0 vector.

If we wanted to account for the bias, it's the same but include a new last row for each matrix, each vector has one more element.

## 2.1 Summary

$$\begin{aligned} Z &= X^+ (\nabla_S \mathcal{L}(S) + X W_Q W_K^\top X^\top) (X^+)^{\top} \\ &= X^+ \nabla_S \mathcal{L}(S) (X^+)^{\top} + X^+ X W_Q W_K^\top X^+ X \end{aligned} \tag{2.1}$$

and

$$\begin{aligned} U_{k'} \Sigma_{k'} V_{k'}^\top &= \text{SVD}_{\text{rank } k'}(Z) \\ \dot{W}_Q^* &= U_{k'} \Sigma_{k'}^{\frac{1}{2}}, \quad \dot{W}_K^* = V_{k'} \Sigma_{k'}^{\frac{1}{2}}. \end{aligned}$$

### 2.1.1 Computing efficiency for $Z$