(a) **Case 1.** $\quad K = Q = X + \mathcal{E}.$

score $\quad K^T Q = (X + \mathcal{E})^T (X + \mathcal{E})$

$$= (X^T + \mathcal{E}^T)(X + \mathcal{E})$$

$$= (X^T X + X^T \mathcal{E} + \mathcal{E}^T X + \mathcal{E}^T \mathcal{E})$$

$$= (X^T X + \mathcal{E}^T \mathcal{E}) + X^T \mathcal{E} + \mathcal{E}^T X$$

terms meaning →

$X^T X$ and $\mathcal{E}^T \mathcal{E}$ → Captures corelations between different features for both feature and positional embedding space.

$\mathcal{E}^T X$ → this matrix calculates projection of each positional embedding vector in the feature space.

$X^T \mathcal{E}$ → this matrix calculates projection of each sample in embedding space.

**Case 2** →. $\quad K = Q = Cat(X, \mathcal{E}).$

$$K = \begin{bmatrix} X & \vdots & \mathcal{E} \end{bmatrix}$$

Sort of like additional Samples in features space.

Cleeeely $K^T Q$ Now will be the Matrix which Calculates covarience matrix but now with 2N samples.