**forward:**

① $\begin{cases} G = KK^T \\ A = I + \text{tril}(\text{diag}(b)G, -1) \end{cases}$

② $B_K = \text{diag}(b)K, \quad B_V = \text{diag}(b)V$

③ $W = A^{-1}B_K, \quad U = A^{-1}B_V$

④ $X = U - WS^T$

⑤ $P = (QK^T) \odot M$

⑥ $O = QS^T + PX$

⑦ $S^+ = S + X^T K$

---

let $\bar{x}$ be $\frac{\partial L}{\partial x}$, upstream $\bar{S}^+, \bar{O}$

⑦ $\bar{S} \mathrel{+}= \bar{S}^+, \quad \bar{X} \mathrel{+}= K(\bar{S}^+)^T, \quad \bar{K} \mathrel{+}= X\bar{S}^+$

⑥ $\bar{Q} \mathrel{+}= \bar{O}S, \quad \bar{S} \mathrel{+}= \bar{O}^T Q$

$\bar{P} \mathrel{+}= \bar{O}X^T, \quad \bar{X} \mathrel{+}= P^T\bar{O}$

⑤ let $T = QK^T \to \bar{T} = \bar{P} \odot M$

$\bar{Q} \mathrel{+}= \bar{T}K, \quad \bar{K} \mathrel{+}= \bar{T}^T Q$

④ $\bar{U} \mathrel{+}= \bar{X}, \quad \bar{W} \mathrel{+}= -\bar{X}S, \quad \bar{S} \mathrel{+}= -\bar{X}^T W$

③ $\bar{B}_K = A^{-T}\bar{W}, \quad \bar{B}_V = A^{-T}\bar{U}$

$\bar{A} \mathrel{+}= -A^{-T}\bar{W}W^T - A^{-T}\bar{U}U^T$

---

$Y = A^{-1}B \to AY = B$

~~$A(dY)+(dA)Y=dB$~~

$A(dY)+(dA)Y=dB$

$dY = -A^{-1}(dA)Y + A^{-1}dB$

$\bar{G} := \frac{dL}{dY}$ By Riesz, $dL = \langle \bar{G}, dY\rangle = \text{tr}(\bar{G}^T dY)$

$dL = \text{tr}\left(\bar{G}^T\left[-A^{-1}(dA)Y + A^{-1}dB\right]\right)$

$= -\text{tr}\left(\bar{G}^T A^{-1}(dA)Y\right) + \text{tr}\left(\bar{G}^T A^{-1}dB\right)$

$\quad\quad \searrow -\text{tr}\left((A^{-T}\bar{G})^T(dA)Y\right)$

$\therefore \frac{dL}{dB} = A^{-T}\bar{G} \quad\quad = -\text{tr}\left(Y(A^{-T}\bar{G})^T(dA)\right)^{[3]}$

$\frac{dL}{dA} = -A^{-T}\bar{G}Y^T \quad = -\text{tr}\left((A^{-T}\bar{G}Y^T)^T(dA)\right)$

---

lower mask $= L_{rs} = \mathbb{1}[r>s]$

$A_{rs} = \delta_{rs} + L_{rs}(b_r G_{rs})$

for $A$ on diagonal or upper is constant so those $\bar{b}$ & $\bar{G} = 0$

$r > s: \quad A_{rs} = b_r G_{rs}$

$G = KK^T \quad\quad dG = dKK^T + KdK^T$

$\bar{G} = \frac{dL}{dG} \quad\quad dL = \langle\bar{G}, dG\rangle$

$\quad\quad\quad = \langle\bar{G}, dKK^T\rangle + \langle\bar{G}, KdK^T\rangle$

$\quad\quad\quad = \langle\bar{G}K, dK\rangle + \langle\bar{G}^T K, dK\rangle$

$\quad\quad\quad = \langle(\bar{G}+\bar{G}^T)K, dK\rangle$

③: trace is cyclic

---

$B_K \in \mathbb{R}^{c \times D} \quad (B_K)_{id} = b_i K_{id}$

$\bar{K}_{id} \mathrel{+}= b_i(\bar{B}_K)_{id}$

$\bar{V}_{id} \mathrel{+}= b_i(\bar{B}_V)_{id}$

$\bar{b}_i \mathrel{+}= \sum_{d=1}^{D}\left((\bar{B}_K)_{id}K_{id} + (\bar{B}_V)_{id}V_{id}\right)$

$\bar{b}_i \mathrel{+}= \sum_{s<i}\bar{A}_{is}G_{is}$

$\bar{G}_{is} = \bar{A}_{is}\frac{\partial A_{is}}{\partial G_{is}} = \begin{cases}\bar{A}_{is}b_i, i>s \\ 0, \text{o/w}\end{cases} = \text{tril}(\text{diag}(b)\bar{A}, -1)$

$\bar{K} \mathrel{+}= (\bar{G} + \bar{G}^T)K$

update $\bar{S}^+ \leftarrow \bar{S}$, then go previous chunk