

Rethinking Attention:

$$\text{Att}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$

Q, K, V have shape $(L \times d)$ (L sequence length)
(d dimension of queries/keys/values)

General time complexity for matrix $A \times B = C$
 $(n \times m) \times (m \times p) = (n \times p)$
complexity $O(nmp)$

If we ignore softmax and say $\text{Att}' = (QK^T)V$ we get
different complexities depending on the matrix multiplication order.

$$\text{Att}' = (QK^T)V$$

 $(L \times d) \times (d \times L) \times (L \times d)$
 $(L \times L) \times (L \times d)$
output shape $(L \times d)$
complexity $O(Ld^2)$

$$\text{Att}' = Q(K^TV)$$

 $(L \times d) \times (d \times L) \times (L \times d)$
 $(L \times d) \times (d \times d)$
output shape $(L \times d)$
complexity $O(Ld^2)$

The latter $O(d^2L)$ is preferred as we can choose $d \ll L$.
 Assume $d \ll L$. $O(L^2d)$ is quadratic complexity wrt.
 sequence length. But this is only possible by removing the softmax.

(can we find an approximation such that $Q'k'^T \approx \text{sm}\left(\frac{Qk^T}{\sqrt{d}}\right)$)

Softmax is a function given vector \vec{z} , of length n , normalizes the elements
 as: $\sigma(z_i) = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$

Kernels: Kernels are functions that are equivalent to the dot product of some feature map ϕ

$\gamma(x, y) = \phi(x)^T \phi(y)$ we would like a function γ that
 avoids computing in ϕ 's high dimensional space.

If we assume A is a kernel matrix with $A_{i,j} = \gamma(q_i, k_j) = e^{(q_i, k_j)^T}$

q_i, k_j are raw vectors in $\mathcal{Q} \times \mathcal{K}$

$$A_{ij} = \gamma(q_i, k_j) = e^{\ell_i k_j^\top} = \phi(q_i)^\top \phi(k_j)$$

Most kernels can be approximated via:

$$\phi(x) = \frac{h(x)}{\sqrt{m}} (f_1(w_1^\top x), \dots, f_l(w_m^\top x), \dots, f_L(w_n^\top x))$$

$h(x)$ & $f_1, \dots, f_L(w_j^\top x)$ are deterministic functions,

w_1, \dots, w_m are random values taken from distribution D , $\phi(x)$

is a vector with shape $(L \times n)$

Gaussian kernel if $h=1$, $\ell=2$, $f_1 = \sin(x)$, $f_2 = \cos(x)$, $D = N(0, I_d)$

so if we draw w from a normal distribution with mean 0 & unit

variance, we obtain the gaussian kernel by using the feature map.

$$\phi_{\text{gauss}}(x) = \frac{1}{\sqrt{m}} (\sin(w_1^\top x), \dots, \sin(w_m^\top x), \dots, \cos(w_1^\top x), \dots, \cos(w_m^\top x))$$

$$\gamma_{\text{guess}} = e^{\left(-\frac{\|x-y\|^2}{z}\right) *}$$

No we want $\gamma_{\text{sm}}(x, y) = e^{(x^T y)}$

$$\gamma_{\text{guess}} = e^{-\frac{\|x-y\|\|x-y\|}{z}} = e^{-\frac{(x^T x - \|x\| \|y\| - \|y\| \|x\| + \|y\|^2)}{z}}$$

for product bound $|X||Y| = X^T Y$

$$\gamma_{\text{guess}} = e^{\frac{(x^T y - \|x\|^2 - \|y\|^2)}{z}}$$

$$\gamma_{\text{sm}}(x, y) = e^{x^T y} = \gamma_{\text{guess}} e^{\frac{\|x\|^2}{z}} e^{\frac{\|y\|^2}{z}}$$

We can reuse the feature map that leads to the gaussian kernel

changing the h function from $h(x) = 1$ to $h(x) = e^{\frac{\|x\|^2}{z}}$

A few issues with this, SM always produces +ve value, the current approximation will give some -ve values, since we are drawing from a gaussian with mean Q_g this will lead to some abnormal behavior.

Attention as fast weight linear db.

Input $x^{(g)} \in \mathbb{R}^d$

$$k^{(i)} = w_k x^{(i)} \quad q^{(i)} = w_q x^{(i)} \quad v^{(i)} = w_v x^{(i)}$$

$$K^{(i)} = [K^{(i-1)}, k^{(i)}]$$

$$V^{(i)} = [V^{(i-1)}, v^{(i)}]$$

$Q K^i V$

but this breaks
commutativity

$$y^{(i)} = V^{(i)} \text{softmax}((K^{(i)})^T q^{(i)})$$

$a \times 6 \neq 6 \times a$

If we ignore the normalisation & softmax

$$(a + (b + c)) = ((a + b) + c)$$

$$y^{(i)} = V^{(i)}((K^{(i)})^T \cdot q^{(i)}) = (V^{(i)}(K^{(i)})^T) \cdot q^{(i)}$$

$$= \left(\sum_{j=1}^i V^{(j)} \otimes K^{(j)} \right) \cdot q^{(i)}$$

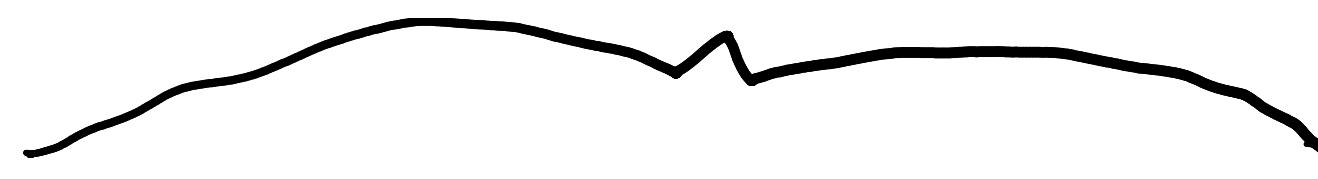
BRA-KET
Notation $|\psi\rangle \langle\psi| = |\psi\rangle^\dagger$ Transpose/conjugate
 $|\phi\rangle \langle\phi| = |\phi\rangle^\dagger$ Inner product
 $\phi^\dagger \circ \psi = \langle\phi|\psi\rangle$ Outer product

$$= \left(\sum_{j=1}^i V^{(j)} \otimes K^{(j)} \right) |q^{(i)}\rangle = W^{(i)} |q^{(i)}\rangle$$

$$W^{(i)} = W^{(i-1)} + |V^{(i)} \times K^{(i)}|$$

$$q^{(i)} = W^{(i)} |q^{(i)}\rangle$$

key / value db entries



linear additive database: $W^{(i)} = |V^{(1)} \times K^{(1)}| + |V^{(2)} \times K^{(2)}| + \dots$

$$= \sum_{j=1}^n |V^{(j)} \times K^{(j)}| \quad (n \text{ entries in DB})$$

Query $|q^{(i)}\rangle$ on database $W^{(i)}$ $W^{(i)} |q^{(i)}\rangle = \left(\sum_j |V^{(j)} \times K^{(j)}| \right) |q^{(i)}\rangle$

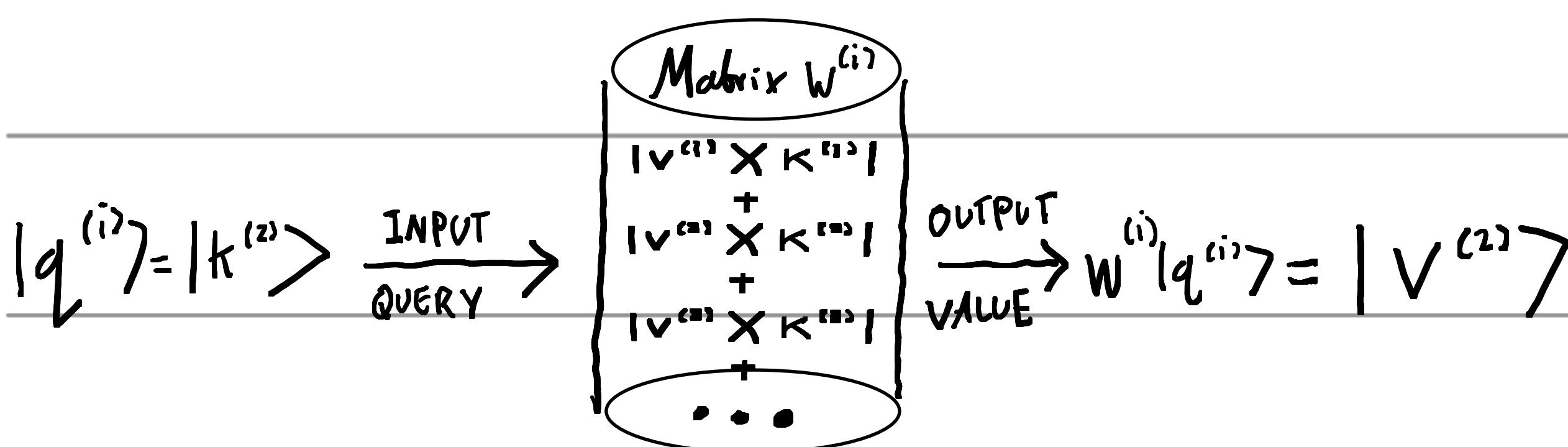
lets say $|q^{(i)}\rangle = |K^{(2)}\rangle$ $W^{(i)} |K^{(2)}\rangle = \left(|V^{(1)} \times K^{(1)}| + |V^{(2)} \times K^{(2)}| + \dots \right) |K^{(2)}\rangle$

$$= |V^{(1)} \times K^{(1)}| + |V^{(2)} \times K^{(2)}| + \dots$$

$$= |V^{(2)}\rangle$$

all zeros

more key / values



The query $|q^{(i)}\rangle$ selects $|V^{(2)}\rangle$ from $W^{(i)}$ given $|K^{(2)}\rangle$, $|K^{(n)}\rangle$ needs

to be orthogonal or you would get a mixture of multiple values.