

PCA ANALYSIS: EFFECTIVE BUFFER DIMENSIONALITY

$$X = \begin{bmatrix} & \\ N & \end{bmatrix} \quad D$$

GRAM MATRIX : the matrix of scalar products

$$\left\{ \begin{array}{l} G_N := X \cdot X^T \\ (N \times D) \quad (D \times N) \end{array} \right.$$

$$G_D := X^T \cdot X \quad (D \times N) \quad (N \times D)$$

$$g_D(i,j) = \text{Cov}(x_i, x_j)$$

$$0 \leq \text{rk}(X) \leq N$$

$$X = \begin{bmatrix} \xrightarrow{\quad} \\ \xrightarrow{\quad} \\ \xrightarrow{\quad} \end{bmatrix}$$

$\text{span}\{\text{rows}(X)\}$ is a subspace of \mathbb{R}^D

The maximum dimension is N .

We want to find its dimension = the rank of matrix X .

STRATEGY: 1) center X : $X = X - \bar{x}_{\text{mean}}$ (-1 degree of freedom row)

2) we apply PCA on the (normalized) Gram matrix $G_D := \frac{1}{N-1} X^T X$. This has max rank $(N-1)$. The number of non-null eigenvalues of $\mathcal{L} = U G_D V^T$ gives the true dimension of X .

Max dim of N points is $N-1$ So?

$$\left\{ \begin{array}{l} \text{max dim}(\vec{x}) = 0 \quad (\text{one point}) \\ \text{max dim}(\vec{x}_1, \vec{x}_2) = 1 \quad (\text{a straight line}) \\ \text{max dim}(\vec{x}_1, \vec{x}_2, \vec{x}_3) = 2 \quad (\text{a plane}) \\ \vdots \end{array} \right.$$

CAVEATS OF THE FINITE SAMPLE

- \mathcal{L} has N eigenvalues \rightarrow discrete spectrum.
- for data $N \rightarrow +\infty$, \mathcal{L} has a continuous spectrum
- for data $N \rightarrow +\infty$, in particular, at one point $N > D$: therefore, the matrix X can reach rank D , which would mean that the true dimensionality of the buffer is really D and not lower

CASE A



for moderate N there seems to exist a subspace $D < N$ where points live

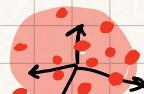


adding more and more points, it stays the same dimension

CASE B



for moderate N there seems to exist a subspace $D < N$ where points live



but it disappears adding more and more points

PCA: an essential RECAP

$$\mathcal{X} = \begin{pmatrix} \parallel \\ \parallel \\ \parallel \end{pmatrix} \quad N \text{ indep. observations of variables } \vec{x} = (x_1, \dots, x_D)$$

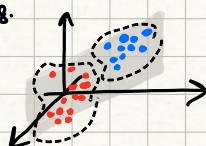
we want to study the covariance between pairs of random variables. Easliy show that this is given by

$$\Sigma := \underset{\text{sample estimate}}{\widehat{\text{Cov}}}(X_i, X_j) = \frac{(\bar{X} - \bar{X})(\bar{X} - \bar{X})^T}{N-1}.$$

EXAMPLE:

features that don't vary or vary little are features we should care about (e.g. if we want to do clustering)

$$\begin{cases} X_1 = \mathcal{N}(0, 1) \\ X_2 = \mathcal{N}(0, 3) + X_1 \\ X_3 = 2X_1 + X_2 \end{cases} \quad \begin{matrix} \text{determined} \\ \text{completely from the} \\ \text{other two} \end{matrix}$$



Σ is the $N \times N$ gram matrix: $\Sigma = \left(\frac{1}{N-1} \right) \mathcal{X}^T \mathcal{X}$
 $(N \times 1) (N \times N) (N \times 1)^T$

Being square and sym, it can be diagonalised:

$\Lambda = Q^T \Sigma Q$. Λ is the covariance matrix of the transformed data $\mathcal{X}' = \mathcal{X} \cdot Q$. In the new basis B ($T_{\mathcal{E}^B} = Q^T$) the random variables are independent because Λ is diagonal.

In the rotated frame, $\sigma^2(x_3') = 0$ (one coordinate is always zero)

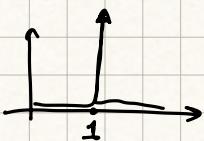
The number of non-null eigenvalues gives the effective dimensionality of the dataset.

EIGENVALUE DENSITY DISTRIBUTION

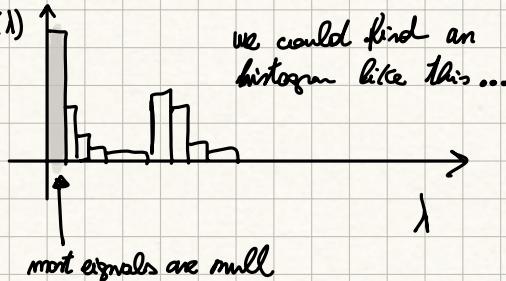
I_N identity matrix

$\lambda = 1$ identity

all N eigenvalues
are 1



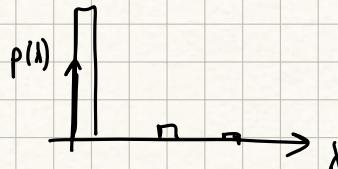
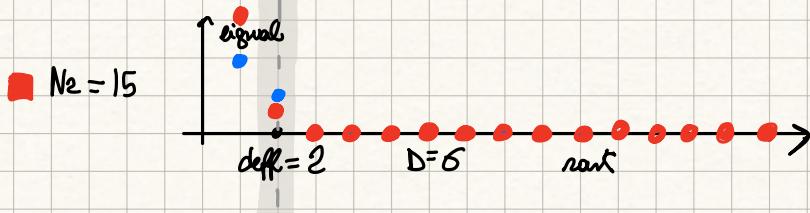
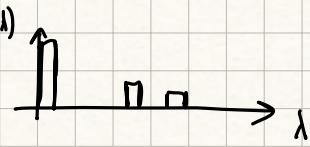
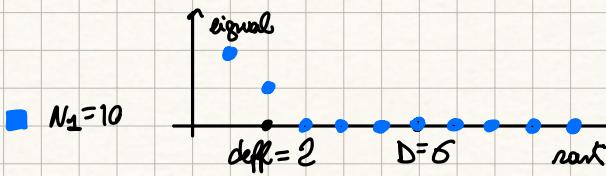
lets estimate $p(\lambda)$
from the observed
empirical frequencies
of each signal



we should ask ourselves:

what would it
become like in the
large- N limit?

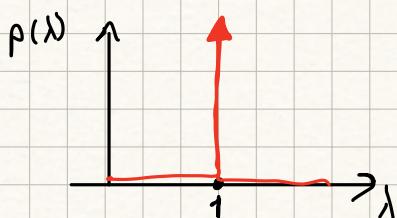
FIG. VS. RANK PLOT IF THE BUFFER REALLY HAS AN EFFECTIVE DIMENSIONALITY $\text{def} < D$



limit $p(\lambda) = \delta(\lambda)$

(all last def are zero)

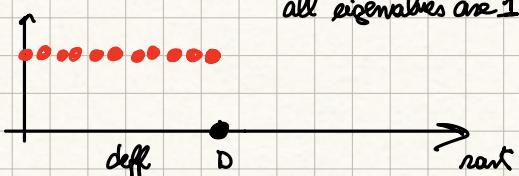
QUESTION : Why check $X^T X$ against I_D ?



$$p(\lambda) = \delta(\lambda - 1)$$

other possibilities ...

$$I_D = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

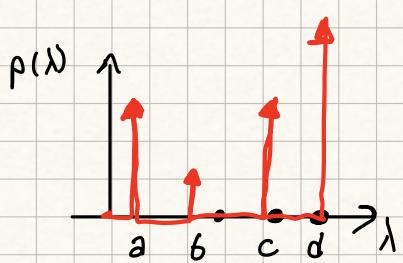


all eigenvalues are 1

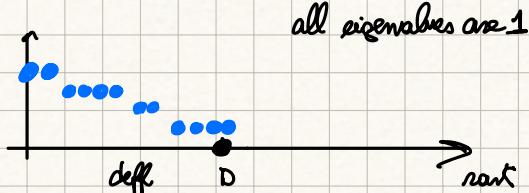
equal variance in all directions \Rightarrow max dimension, D



ball of noise (with FINITE variance! - doesn't span over all the available space)



$$I_D = \begin{pmatrix} a & 0 & 0 & 0 \\ 0 & b & 0 & 0 \\ 0 & 0 & c & 0 \\ 0 & 0 & 0 & d \end{pmatrix}$$



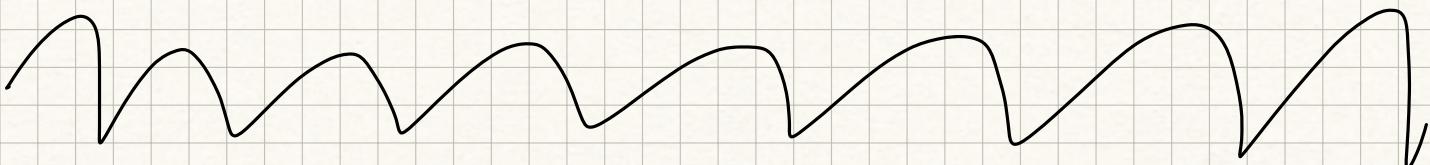
all eigenvalues are 1

equal variance in all directions \Rightarrow max dimension, D

$$\begin{aligned} p(\lambda) = & \alpha \cdot \delta(\lambda - a) + \\ & \beta \cdot \delta(\lambda - b) + \\ & \gamma \cdot \delta(\lambda - c) + \\ & \zeta \cdot \delta(\lambda - d) \end{aligned}$$

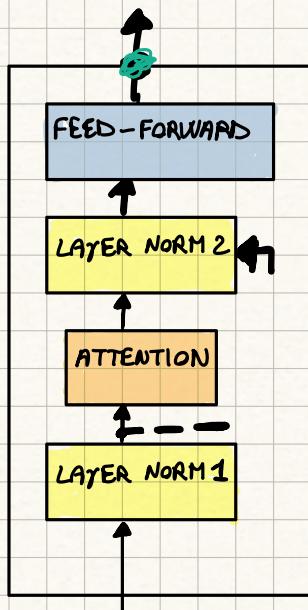
I think all we said doesn't really apply!

we need to look at $X^T X$ for dimensionality!
not the other!



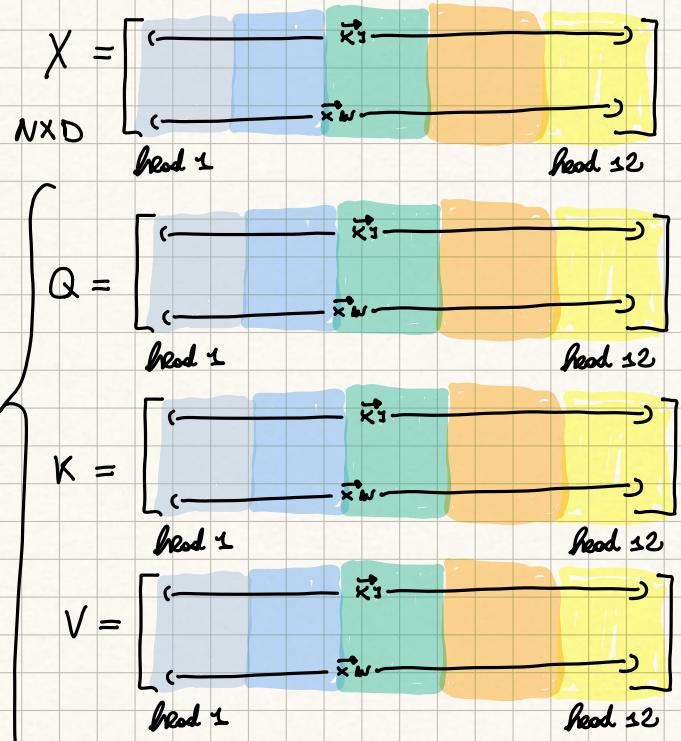
EQUATIONS: WHAT IS HAPPENING EXACTLY, STEP BY STEP

THE ATTENTION LAYER



$$\begin{array}{ll} \text{C-proj-bias} & D \\ \text{C-proj-weight} & D \times D \\ \text{C-attn-bias} & 3D \\ \text{C-attn-weight} & D \times 3D \end{array}$$

$$n_{\text{HEAD}} = 12, \quad d_h = \frac{D}{n_{\text{HEAD}}} = 64$$



STEP 1 : create queries, values and keys

$$(C\text{-attn-weight}, C\text{-attn-bias})$$

$$D \times 3D \quad 3D$$

$$[Q, K, V] = X \cdot C\text{-attn-weight} + \begin{bmatrix} C\text{-attn-bias} \\ C\text{-attn-bias} \\ C\text{-attn-bias} \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} C\text{-attn-bias} \\ C\text{-attn-weight} \end{bmatrix}$$

stacked horizontally

equivalent

STEP 2 : Split into heads

NOTICE : each chunk still carries inside the information coming from the whole buffer X.

$$Q = \begin{bmatrix} \cdots & \xrightarrow{\vec{K}_1} & \cdots & \xrightarrow{\vec{K}_{12}} & \cdots \end{bmatrix} \rightarrow \begin{bmatrix} \cdots \\ \cdots \end{bmatrix}, \begin{bmatrix} \cdots \\ \cdots \end{bmatrix}, \begin{bmatrix} \cdots \\ \cdots \end{bmatrix}$$

$\nabla Q_h = Q_h(X)$ and not $Q_h(X_h)$. It is a non linear operation.

STEP 3: compute attention weights and stack heads together

$$\text{head}_i = e^{\left(\frac{Q_h \cdot K_h^T}{\sqrt{d_k}} \right)} \cdot V_h \rightarrow [\text{head}_1, \text{head}_2, \dots, \text{head}_{12}]$$

$$N \times d_k \quad N \times N \quad N \times d_k \quad N \times (d_k \times 12) = N \times D$$

STEP 4: "inject back" ?

$$Y = [1, \text{head}_1, \text{head}_2, \dots, \text{head}_{12}] \cdot \begin{bmatrix} C\text{-proj-bias} \\ C\text{-proj-weight} \\ (D+1) \times D \end{bmatrix}$$

PROPOSALS OF ANALYSIS

1) QUESTION

Is $Q = V \Lambda Q X$ a projection to a lower dim. space?

$$\text{like in... } Y = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ 0 \end{pmatrix} \text{ or isn't it?}$$



To understand, we should check if the matrices Q, K and V are FULL RANK or not.

MY GUESS: They are NOT full rank.

2) $Q \cdot K^T = X V \Lambda V \Lambda K^T X$. How different from a symmetric matrix? How different from the identity?

If it was symmetric, the machinery would be over-complicated, because all would reduce to first transform $X := XA$ and then take the $n \times n$ Gram matrix: $(X')(X')^T = (XA)(XA)^T = XAA^TX = XMX$ (Choleski decomposition)

Even more, if it was $M \cong I$, then attention weights would be nothing but scalar products of X on itself.

{ GUESS Since $(QK^T)^T = (XM^TX^T)^T = XM^TX'$, QK^T is symmetric $\Rightarrow M$ is symmetric.

{ Make some prompts and show the attention weights: they usually are very asymmetric. Then also M is asymmetric, otherwise we would assume QK^T symmetric ALWAYS.

SPOILER: Indeed, M is ASYMMETRIC and differs from I . Therefore, we know at least that the attention stage is doing something (non-trivial), even if we don't yet understand WHAT.

3) STEP-BY-STEP DIMENSION ANALYSIS (?)

• $[Q, K, V] = X \cdot \text{C-orth-weight dimension of } Q, K, V ?$

• $\text{head}_h = \underbrace{\frac{1}{Z} \exp\left(\frac{Q_h \cdot K_h^T}{\sqrt{d_k}}\right)}_{\text{normalization to 1}} V_h = ?$ not decomposable more than this

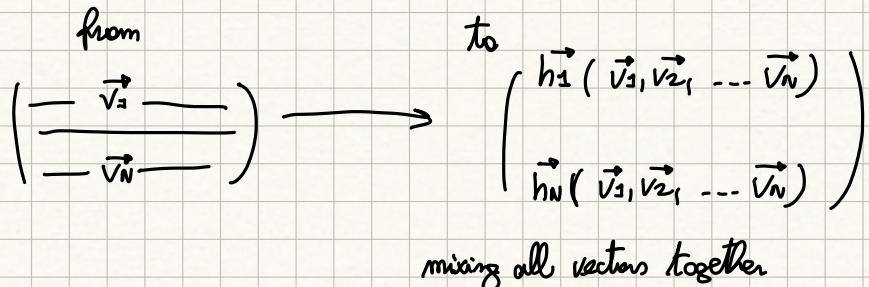
• $Y = [\quad] \cdot V_h$

QUESTION KEY IDEA
shift the \vec{x}_i vector towards its most aligned value \vec{v}_s

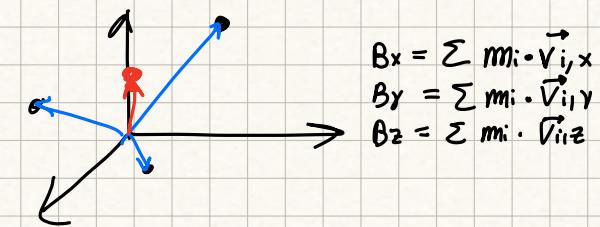
QK^T alignment and magnitude both matter
(angle)

$$V \rightarrow H = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \cdot V \quad h_i = \sum_j \text{weight}_{ij} \vec{v}_j$$

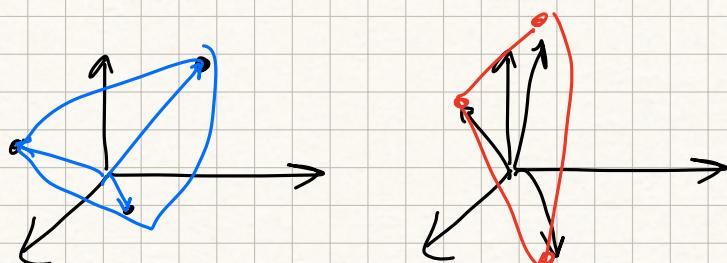
where weights sum up to 1



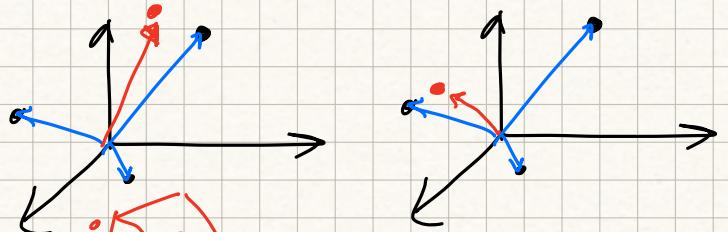
BARYCENTER



each h_i is a differently weighted barycenter of the vectors $\{\vec{v}_i\}$



set of values gets transformed in set of barycenters



QUESTION: what is the volume spanned by the barycenters?
smaller or larger?

ASYMMETRY OF THE WEIGHT MATRIX and DIRECTION OF MOVEMENTS

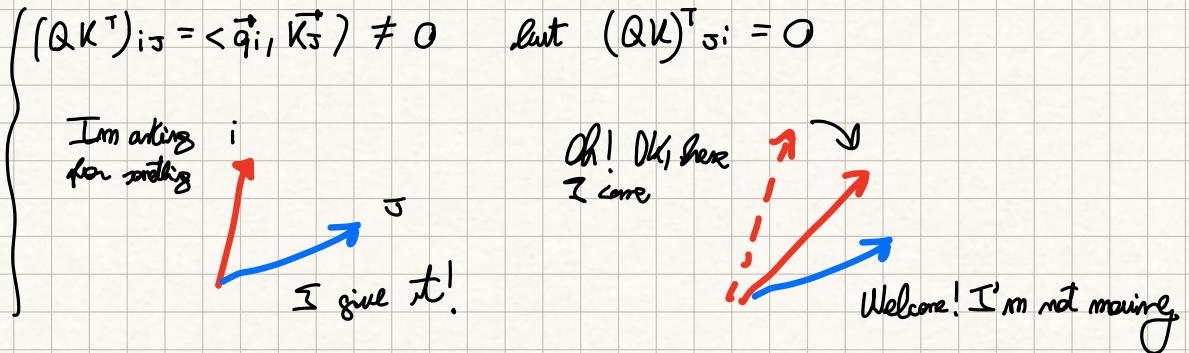
$$\vec{h}_i = \sum_j \text{weight}_{ij} \cdot \vec{v}_j$$

Notice that the fact that the matrix QK^T is asymmetric has a strong consequence:

The magnitude of weight_{ij} depends on how well the query of \vec{x}_i , \vec{q}_i matched against the key of \vec{x}_j (what it had to offer), not viceversa!

$$(QK^T)_{ij} = \langle \vec{q}_i, \vec{k}_j \rangle$$

EXAMPLE $(QK^T)_{ij} = \langle \vec{q}_i, \vec{k}_j \rangle \neq 0$ but $(QK)^T_{ji} = 0$



Vectors move towards those vectors that have something to offer.

In MASKED attention, query of i cannot match key of j if j comes after, because i doesn't see j .

If we trained the model without restricting what it can see, it would not be able to auto-complete text when used in evaluation mode.

$$QK^T = \underset{i}{\left[\begin{array}{c} \text{L} \\ \text{O} \\ \text{J} \end{array} \right]}$$

\vec{x}_i cannot ask anything to the vectors that come AFTER it.

I swam across the river
learning

I swam across the ...
auto-completion task

Say if I understood correctly or not. Causal masked attention means that the query of token i -th is not allowed to match any key of token j -th if j comes after i . If the mask was not set during training, we would get a model that is not able to perform an auto-completion task, since it didn't learn to complete sentences

Yes — you understood exactly right.

Let's restate it in your words with some clarifications:

Your understanding:

Causal masked attention means that the **query** of the i -th token is not allowed to match any **key from tokens** $j > i$.

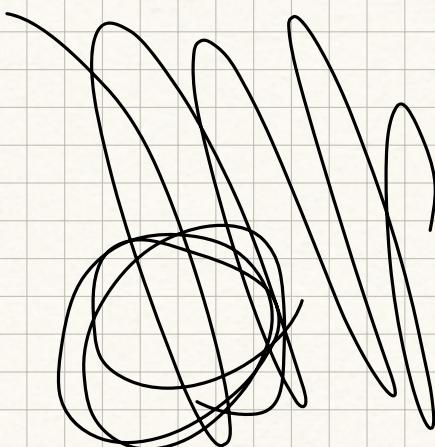
Correct: It **blocks access to future tokens** during training by masking their keys — ensuring token i can only attend to tokens $\leq i$.

Also correct:

If the mask was not set during training, the model would **see the future**, and thus fail to learn auto-completion.

Right: Without the mask, the model would "cheat" — using future tokens as context — and **wouldn't learn to predict the next token from only past information**.

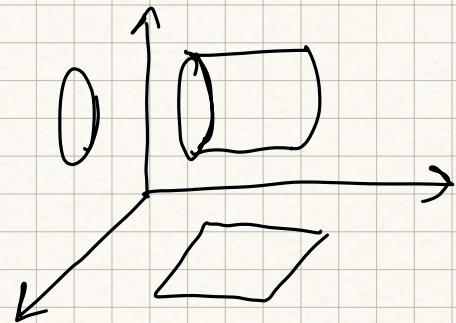
So during inference (when future tokens aren't available), it would fail, because it never learned to predict step-by-step.



But then, do we also need the mask during inference?

MULTI-HEAD ATTENTION

Reads partition elements in different lower-dimensional subspaces



$$Q_h = X \cdot W_{Q,h}$$

(N, D) (D, d_h)

max rank is d_h

it's like looking at an object FROM DIFFERENT POINTS OF VIEW

Like in technical drawings at school

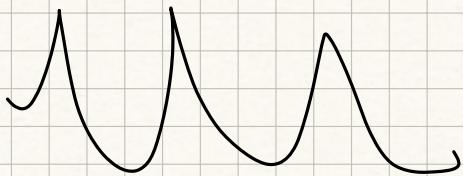
What is

$X \cdot A$ doing on X ?

$N \times D$ $D \times D$

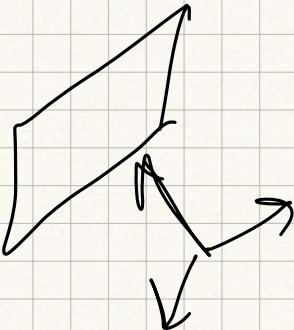
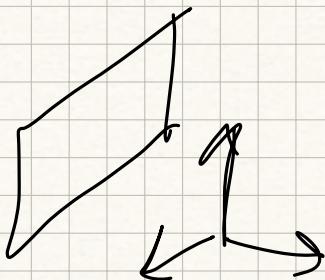
- When A is full rank, a change of basis

• when A is not, a projection.



here projections are not orthogonal
because there is first a change of basis (or is it a projection?)

⚠ CHECK RANK of W_Q, W_K, W_V

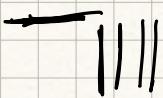


if points are lying on a plane, and you change the basis, they are still lying on a plane.

so

$$Q_n = \frac{\partial}{\partial x} \left[X \cdot \frac{\partial Q_n}{\partial x} \right]$$

$$X \cdot [W_{k_1}, W_{k_2}, \dots, W_{k_n}]$$



$$(Q_1) \quad (Q_2) \quad (Q_n)$$

$$XW_{k_1}$$

$$XW_{k_2}$$

$$XW_{k_n}$$

$$XW_k \cdot XW_k^T = [W_{k_1}, W_{k_2}, \dots, W_{k_n}] \circ \begin{bmatrix} XW_{k_1}^T \\ \vdots \\ XW_{k_n}^T \end{bmatrix} = \sum (XW_{k_i} \cdot XW_{k_i})^T$$

osservare che è possibile
e uguale a proiettare (con matrice
matrice) direttamente.

