

Attention - based

Clustering

quantization?

Joint work with Pierre Marion & Rodrigo Mauel-Soto.

Context: Transformer architecture



Existing literature:

- Mainly in supervised learning (see the work of Peter Battat's team)
- Via in-context learning
- Other work with Pierre Marion, Raphaël Berthier and Gérard Biard where we defined a new statistical task called "single location regression" in which only one token is relevant to form the prediction, but its position is random (and therefore prompt-dependent) . . .

has some link with the ongoing work of Alberto Betti last week.

Motivations:

- Understand the statistical abilities of attention layers in unsupervised scenarios
- Why? Powerful embedding performed by attention layers -
- Existing work by He et al (2025) showing that transformer

can mimic the EM algorithm . But their analysis required the knowledge of the labels during training.

- "Clustering" phenomenon exhibited by the team of detroit, Geshkovski, Polyanskiy & Rigollet. But this is a completely \neq setting
 - ✗ they stack an infinite number of layers and study the continuous dynamics where time = depth
 - ✗ no training of the parameters

THIS TALK: Focus on a single attention layer fully trained in an unsupervised way.

SETTING : MODEL-BASED CLUSTERING

- Input sequence of L tokens where each token $X_l \in \mathbb{R}^d$. $\mathbf{X} = \begin{pmatrix} X_1^T \\ \vdots \\ X_L^T \end{pmatrix} \in \mathbb{R}^{L \times d}$

- Assumption : $X_l \stackrel{iid}{\sim} \frac{1}{2} \mathcal{N}(\mu_0^\star, \sigma^2 I) + \frac{1}{2} \mathcal{N}(\mu_1^\star, \sigma^2 I)$ (P_{T2})

Mixture model with balanced components
centroids as unit vectors

$$\pi_0 = \pi_1 = \frac{1}{2}$$

$$\mu_0^\star, \mu_1^\star \in \mathbb{S}^{d-1}.$$

$$\|\mu_0^\star\|_2 = \|\mu_1^\star\|_2 = 1$$

- Orthogonality assumption

$$\mu_0^\perp \mu_i^\perp \text{ or } \langle \mu_0^\perp, \mu_i^\perp \rangle = 0$$

- For each token x_e , there exists an associated latent variable z_e such that $z_e \sim \mathcal{B}(1/2)$ encodes the corresponding cluster of x_e .

- Degenerate case: $\sigma^2 = 0$ $x_e \sim \frac{1}{2} \delta_{\mu_0^\perp} + \frac{1}{2} \delta_{\mu_1^\perp}$.

ATTENTION-BASED PREDICTOR / EMBEDDING

Attention head: $h^{\text{soft}, \mu}(x) = \text{softmax}_2(x Q K^T x^T) x V$

$$K, Q, V \in \mathbb{R}^{d_{\text{inp}}}$$

Step 1: Take $V = \text{Idl.}$ (no training for V)

Step 2: $K = Q = \mu$ a column matrix $\mu^T \in \mathbb{R}^{1 \times d_{\text{inp}}}$

$$h^{\text{soft}, \mu}(x) = \text{softmax}_2(x \mu \mu^T x^T) x$$

The l -th output vector is therefore given by

$$h^{\text{soft}, \mu}(x)_l = \sum_{k=1}^L \underbrace{\text{softmax}_2(x_k^T \mu \mu^T x_k)}_{\text{we aggregate the } x_k's \text{ when } x_k \text{ and } x_l \text{ are aligned with } \mu.} x_k$$

= we aggregate the x_k 's when x_k and x_l are simultaneously aligned with μ .

Step 3: Getting rid of the softmax

$$H^{\text{lin}, \mu}(\mathbf{x})_l = \frac{2}{L} \sum_{k=1}^L (\lambda \mathbf{x}_k^\top \mu \mu^\top \mathbf{x}_k) \mathbf{x}_k$$

Fix $\mu = \mu^*$. When \mathbf{x}_l and \mathbf{x}_k belong to the cluster O , i.e., such that $z_l = z_k = O$, then the vectors \mathbf{x}_l and \mathbf{x}_k are likely to be aligned with μ^*

$$\|(\mathbf{x}_l^\top \mu^* \mu^* \mathbf{x}_k) \mathbf{x}_k - \mathbf{x}_l\|$$

Conversely, if \mathbf{x}_l and \mathbf{x}_k are associated with \neq clusters, i.e., $z_l \neq z_k$ then $\|(\mathbf{x}_l^\top \mu^* \mu^* \mathbf{x}_k) \mathbf{x}_k - \mathbf{0}\|$

Remark: If $\mu = \mu^*$
 $\mathbf{x}_l \in \text{cluster } O \quad (z_l = O)$
 then $\sum_k \lambda (\mathbf{x}_k^\top \mu^* \mu^* \mathbf{x}_k) \mathbf{x}_k$ aggregates
 the \mathbf{x}_k 's from the same cluster, whose
 expected number is $1/2$ \rightarrow renormalization factor

$H^{\text{lin}, \mu^*}(\mathbf{x})_l = \text{"empirical mean" of the takers}$
 belonging to the same cluster.

Assuming that the # clusters is known, it seems natural to consider an attention-based layer composed of α attention heads

$$T^{lin, \mu_0, \mu_1}(x) = h^{lin, \mu_0}(x) + h^{lin, \mu_1}(x)$$

Metric loss As no label is available, one may consider

$$\begin{aligned} \mathcal{L}(T) &= \frac{1}{L} \sum_{l=1}^L \mathbb{E} \|x_l - T(x)_l\|_2^2 \\ &= \mathbb{E} \|x_q - T(x)_q\|_2^2 \end{aligned}$$

- Note that if T was able to return for each token x_l the corresponding centroid μ_l^* , the risk would correspond to a **QUANTIZATION ERROR**.
- It also looks like the kind of loss used during self-supervised learning but \neq

Training of the attention. layer

$$\text{Set } \mathcal{R}(\mu_0, \mu_1) = \mathcal{L}(T^{lin, \mu_0, \mu_1})$$

PGD iterates
on $(\tilde{D}^{d-1})^2$

$$\mu_0^{k+1} = \frac{\mu_0^k - \tau(I - \mu_0^k (\mu_0^k)^T) \nabla_{\mu_0} \mathcal{R}(\mu_0^k, \mu_1^k)}{\|\mu_0^k - \tau(I - \mu_0^k (\mu_0^k)^T) \nabla_{\mu_0} \mathcal{R}(\mu_0^k, \mu_1^k)\|_2}$$

$$\mu_1^{k+1} = \frac{\mu_1^k - \gamma(I - \mu_1^k(\mu_1^k)^T)^T \nabla_{\mu_1} R(\mu_0^k, \mu_1^k)}{\|\mu_1^k - \gamma(I - \mu_1^k(\mu_1^k)^T)^T \nabla_{\mu_1} R(\mu_0^k, \mu_1^k)\|}$$

TRAINING DYNAMICS IN THE DEGENERATE CASE.

- Reparameterization $\kappa_0 = \langle \mu_0, \mu_0^* \rangle$ $\kappa_1 = \langle \mu_1, \mu_1^* \rangle$

- $\eta_0 = \langle \mu_1, \mu_0^* \rangle$ $\eta_1 = \langle \mu_0, \mu_1^* \rangle$

- Expression of the risk as a polynomial of $\kappa_0, \kappa_1, \eta_0, \eta_1$. (order 4)

- Characterization of global minima

When $\lambda = \frac{L+1}{L+3}$, argmin of the risk characterized by

$$\kappa_0^2 + \eta_0^2 = 1 \quad \kappa_1^2 + \eta_1^2 = 1 \quad \kappa_0 \eta_1 + \kappa_1 \eta_0 = 0 .$$

* computation can be conducted to any value of λ .

- Convergence analysis

Thm: Under the degenerate mixture model, take $\lambda \in [0; \frac{L+1}{L+3}]$,

$\exists \bar{\gamma} > 0$ such that for any stepsize $0 < \gamma < \bar{\gamma}$

for a generic init^o $(\mu_0^o, \mu_1^o) \in \mathcal{O} := \{(\mu_0, \mu_1) \in (\mathbb{S}^{d-1})^2 : \langle \mu_0^*, \mu_1 \rangle = 0, \langle \mu_1^*, \mu_0 \rangle = 0, \langle \mu_0, \mu_1 \rangle = 0\}$

$$(\mu_0^k, \mu_1^k) \xrightarrow{k \rightarrow \infty} (\pm \mu_0^*, \pm \mu_1^*)$$

- G Despite the non-convexity of the pb, the key/query matrices trained via PGD align with the centroids of the underlying Dirac mixture.
- G Initialization on \tilde{S} impractical !!
- G In practice

(a) CV to the centroids observed for init⁰ on \tilde{S}

(b) Stagnation $\longrightarrow (\mathbb{S}^{d-1})^2$

(c) To mitigate this effect, we introduce a regularization

$$\pi(\mu_0, \mu_1) = \mathbb{E}[\langle \mu_0, x_i \rangle^2 \langle \mu_1, x_i \rangle^2]$$

so that training is done via PGD for

$$\min_{\mu_0, \mu_1, \mathbb{E}[(\mathbb{S}^{d-1})^2]} \mathcal{R}(\mu_0, \mu_1) + \rho \pi(\mu_0, \mu_1)$$

ENFORCES THE ORTHOGONALITY

AND HELPS THE ATTENTION PARAMETERS TO
CONVERGE TO THE CENTROIDS WHEN
INITIALIZED IN THE SPHERE

SKIPPED .

TRAINING DYNAMICS in THE NON-DEGENERATE CASE

- Setting $\left\{ \begin{array}{l} X \sim \frac{1}{2} N(\mu_0^*, \sigma^2 I) + \frac{1}{2} N(\mu_1^*, \sigma^2 I) \\ \mu_0^*, \mu_1^* \in \mathbb{S}^{d-1} \text{ and } \langle \mu_0^*, \mu_1^* \rangle = 0 \end{array} \right.$

- Reparameterization $K_0 = \langle \mu_0, \mu_0^* \rangle$ $K_1 = \langle \mu_1, \mu_1^* \rangle$
 $\eta_0 = \langle \mu_1, \mu_0^* \rangle$ $\eta_1 = \langle \mu_0, \mu_1^* \rangle$
 $\xi = \langle \mu_1, \mu_0 \rangle$

- Risk = polynomial of order 4 of these quantities

- CV analysis

Thm: Under the GMM, take $\gamma \in]0; \pi^*(\gamma, L)]$

$\exists \bar{\gamma}$ s.t. for any step size $0 < \gamma < \bar{\gamma}$, for a generic init $(\mu_0^*, \mu_1^*) \in \partial \gamma := \left\{ (\mu_0, \mu_1) \in (\mathbb{S}^{d-1})^2 : \begin{array}{l} \langle \mu_0^*, \mu_1 \rangle = 0 \\ \langle \mu_1^*, \mu_0 \rangle = 0 \\ \langle \mu_1, \mu_0 \rangle = 0 \end{array} \right\}$

$$(\mu_0^k, \mu_1^k) \xrightarrow{k \rightarrow \infty} (\pm \mu_0^*, \pm \mu_1^*)$$

value of γ s.t.
 $(\pm 1, \pm 1)$ are global min.

- ✓ Attention layers can uncover and encode latent structure of the input distribution in a fully unsupervised setting
- ✓ Interpretability of the attention parameters.

- G Despite the non-convexity of the pb, the key / query vectors trained via PGD align with the underlying centroids of the mixture .
- G Initialization on \mathcal{C} impractical !!
- G In practice

- ① CV to the centroids observed for init⁰ on \mathcal{C}
- ② ~~Stagnation~~ $\longrightarrow (\mathbb{S}^{d-1})^2$

- ③ To mitigate this effect , we introduce a regularization

$$r(\mu_0, \mu_1) = \mathbb{E} [\langle \mu_0, x_i \rangle^2 \langle \mu_1, x_i \rangle^2]$$

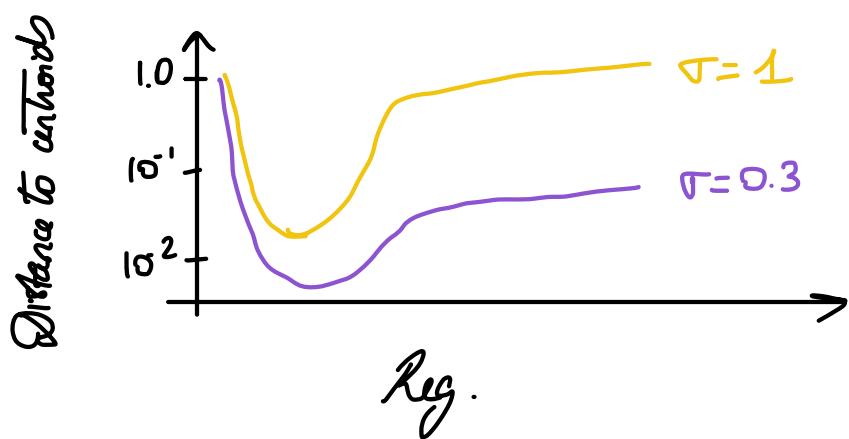
so that training is done via PGD for

$$\min_{\mu_0, \mu_1, \in (\mathbb{S}^{d-1})^2} \mathcal{R}(\mu_0, \mu_1) + \rho r(\mu_0, \mu_1)$$

ENFORCES THE ORTHOGONALITY

AND HELPS THE ATTENTION PARAMETERS TO
CONVERGE TO THE CENTROIDS WHEN
INITIALIZED IN THE SPHERE

G XP



ATTENTION-BASED LAYERS AS APPROXIMATE QUANTIZERS

- Motivation: understanding the statistical properties of an attention layer which parameters have converged to the true centroids.
- Consider the "optimal" quantizer $T^*(x) = \mu_{z_e}^*$.

- The risk of the optimal quantizer

$$\begin{aligned}\mathcal{L}(T^*) &= \mathbb{E}[\|x - \mu_{z_e}^*\|^2] \\ &= \sigma^2.\end{aligned}$$

Lemma: Under the GMN, it holds that

$$\mathbb{E}[T^{lin, \mu_0^*, \mu_1^*}(x)_1 | z_1=c] = \mu_c^* \frac{\lambda}{L} \left[(L+1) + 2(L+3)\sigma^2 \right]$$

Choosing $\lambda = \frac{L}{(L+1) + 2(L+3)\sigma^2}$ gives $\mathbb{E}[T^{lin, \mu_0^*, \mu_1^*}(x)_1 | z_1=c] = \mu_c^*$. (*)

- ✓ The l -th output aligns on average with the centroid of the cluster to which the l -th token belongs.

Prop: Under the GMM, fix $\lambda = \frac{1+4\sigma^2 + 4\sigma^4}{1+6\sigma^2 + 12\sigma^4 + 8\sigma^6}$,
then

$$\mathcal{L}(T^{lin, \mu_0^*, \mu_1^*}) \underset{L \rightarrow \infty}{\sim} (d-d) \sigma^2$$

striking obs^o

$$\frac{\mathcal{L}(T^{lin, \mu_0^*, \mu_1^*})}{\mathcal{L}(T^*)} \underset{L \rightarrow \infty}{\sim} 1 - \frac{2}{d}.$$

!! For long input sequences, attention layer can achieve a better risk than the optimal quantizer

!! But the comparison is not entirely fair:

① $T^{lin, \mu_0^*, \mu_1^*}$ is not a quantizer per se

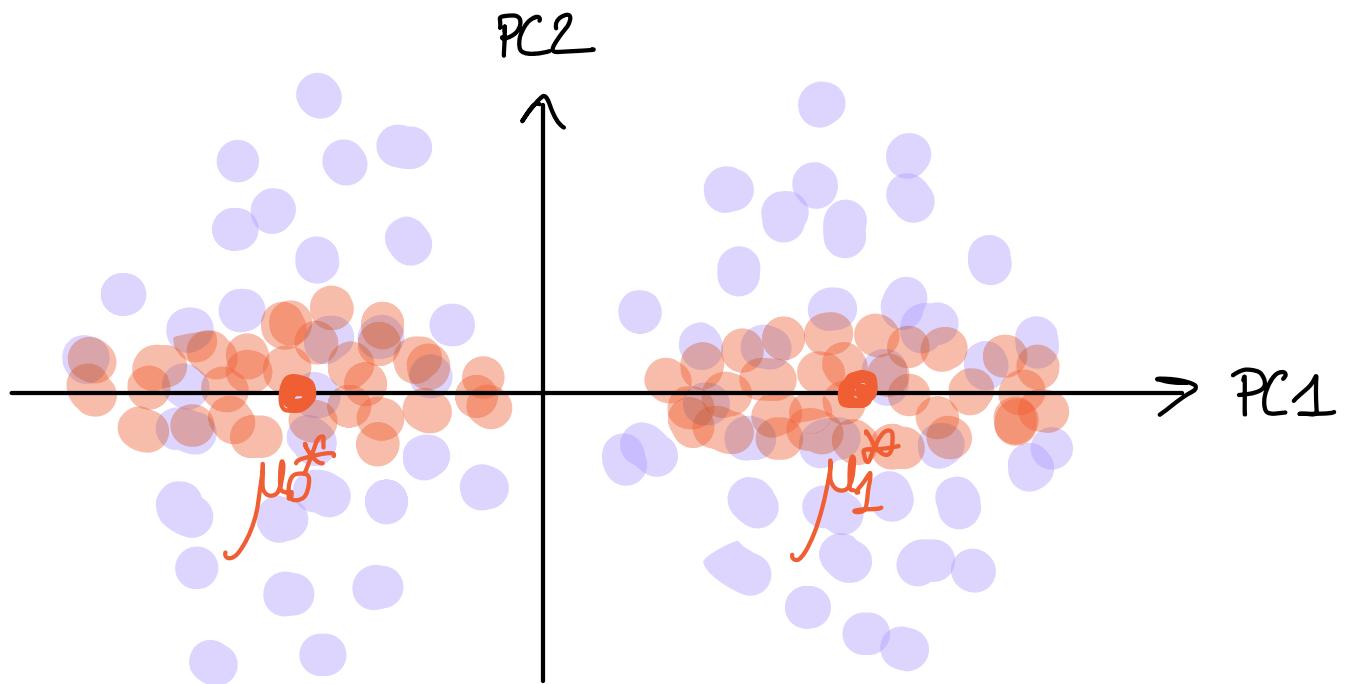
② the optimal quantizer T^* relies on a single input token

to predict the associated centroid (being oracle w.r.t centroids)

(vs.) the oracle attention layer $T^{lin, \mu_0^*, \mu_1^*}$ uses the info of the underlying centroids but benefits from a growing sequence of random var., all drawn from the same Gaussian mixture. \rightarrow VARIANCE REDUCTION MECHANISM.

⊗ with the same choice of λ : $V\text{ar}[T^{lin, \mu_0^*, \mu_1^*}(x)_1 | z_1] = 2\sigma^2$
 instead of d

Illustration of the variance reduction mechanism.



EXPLORATION : IN-CONTEXT CLUSTERING / QUANTIZATION

- Up to now, with 2 heads, the transformed 1st token is

$$T^{\text{lin}, \mu_0, \mu_1}(x) = \frac{2}{L} \sum_{l=1}^L \lambda (x_l^\top \mu_0 \mu_0^\top x_l) x_l + \lambda (x_l^\top \mu_1 \mu_1^\top x_l) x_l.$$

$$= \frac{2}{L} \sum_{l=1}^L \lambda \left[x_l^\top (\mu_0 \mu_0^\top + \mu_1 \mu_1^\top) x_l \right] x_l$$

when $\mu_0 \perp \mu_1$ this is a rank-2 matrix

- Remark : we just showed that we could train rank-2 key / query head by leveraging the non-convex optimization of 2 simple raw-structured heads -

- Challenge : IN-CONTEXT FRAMEWORK

Imagine now that for each input sequence, the centroids μ_0^* and μ_1^* are random.

If they are random but still aligned with particular directions μ_0^{**} and μ_1^{**} , then $T^{\text{lin}, \mu_0, \mu_1}$ should remain a good "quantized" / "embedding" candidate.

Imagine now that for n sequences of L tokens, the centroids are randomly drawn on the sphere:

$$\mu_0^* \sim \text{Un}(S^{d-1})$$

$\mu_1^* | \mu_0^*$ distributed on $S^{d-1} \cap (\mu_0^*)^\perp$.

 $T^{\text{lin}, \mu_0, \mu_1}$ will struggle to capture the changing structure of the data because of the use of only 2 param. $\mu_0, \mu_1 \in S^{d-1}$

Idea: Increase the degrees of freedom of the attention layer
take $\mu_1, \dots, \mu_d \in S^{d-1}$ of unit-norm
mutually orthogonal.

In such a case, the attention layer has no more parameters to learn !!

$$T^{\text{ctr}}(\mathbf{x})_l = \sum_{c=1}^d T^{\text{lin}, \mu_c}(\mathbf{x})_l = \frac{2\lambda}{L} \sum_{c=1}^d \sum_{k=1}^L \left(\mathbf{x}_k^T \mu_c \mu_c^T \mathbf{x}_k \right) \mathbf{x}_k$$

$$= \frac{2\lambda}{L} \sum_{k=1}^L \mathbf{x}_k^T \left(\sum_{c=1}^d \mu_c \mu_c^T \right) \mathbf{x}_k$$

$$= \frac{2\lambda}{L} \sum_{k=1}^L \mathbf{x}_k^T \mathbf{x}_k \xrightarrow{\text{Id}}$$

$$= \frac{2\lambda}{L} \overbrace{\left[\sum_{k=1}^L \mathbf{x}_k \mathbf{x}_k^T \right]}^{\hat{\Sigma}} \mathbf{x}_k$$

To simplify, we could remove the "autocorrelated" term

$$T^{cte}(x)_l = \frac{2\lambda}{L} \sum_{k \neq l} X_l^T X_k X_k.$$

- What can be said about this embedding?

$$\mathbb{E}[T^{cte}(x)_l \mid \mu_1^*, \mu_0^*, z_l = c] = \frac{2\lambda}{L} C_{d, \sigma^2} \mu_c^* \xrightarrow{1 + (d+2)\sigma^2 + (L+1) \left(\frac{1}{2} + \sigma^2 \right)}$$

- ✓ Aligned with μ_c^* in expectation
- ✓ One can choose λ to get "unbiased" encoding with centroids

$$\text{Var}[T(x)_l \mid \mu_1^*, \mu_0^*, z_l = c] \underset{L \rightarrow \infty}{\sim} 2\lambda^2 \sigma^2 (1 + 4\sigma^2 + 2\sigma^4)$$

Choosing $\lambda = 1/(1+2\sigma^2)$ (unbiased encoding when $L \rightarrow \infty$)

$$2\sigma^2 \frac{1 + 4\sigma^2 + 2\sigma^4}{(1 + 2\sigma^2)^2} \leq \sigma^2 d$$

With an appropriate choice of the temperature τ , one can show that

$$\frac{1+2\tau^2}{1+6\tau^2+12\tau^4+4d\tau^6} \leq 1.$$

$$\lim_{L \rightarrow \infty} \frac{\mathcal{L}(T^{dx})}{\mathcal{L}(T^\infty)} = \left(1 - \frac{2}{d}\right) C_{T,d} \stackrel{1 \leq 1}{\leq} 1 - \frac{2}{d}.$$

Better quantization risk than the optimal quantizer and the improvement factor is now dimension-dependent.

Identification of a similar phenomenon

- in a more complex setting (in context)
- with a simpler architecture (no trained parameters)

This simple (non-parametric) encoder manages to capture structure of the data