

# Softmax as Linear Attention in the Large-Prompt Regime: a Measure-based Perspective

with C. Boyer

## Summary:

- Measure-based view of self attention unifying finite and infinite prompts inputs
- Concentration of the finite prompt regime to the infinite one as  $L \rightarrow \infty$
- Derive new convergence guarantees for in-context linear regression with softmax attention

## I - Measure based view of attention

### Classic view

An attention layer takes as inputs a prompt of  $d$ -dimensional tokens  $(z_1, \dots, z_L) \in \mathbb{R}^{L \times d}$  and a query token  $z \in \mathbb{R}^d$  and outputs

$$\frac{\sum_{i=1}^L \exp(\langle K z_i, Q z \rangle) V z_i}{\sum_{i=1}^L \exp(\langle K z_i, Q z \rangle)}$$

} this is a softmax operator  
combined with matrix multiplication

where  $K, Q, V$  are the model parameters (key, query and value matrices)

### Measure based view

Output is independent of token order, so that we can see the prompt as a measure  $\mu \in \mathcal{M}(\mathbb{R}^d)$  and an attention layer parametrized by  $K, Q, V$  is then a function:

$$T^{K, Q, V}: \mathcal{M}(\mathbb{R}^d) \times \mathbb{R}^d \longrightarrow \mathbb{R}^d$$
$$(\mu, z) \longmapsto T^{K, Q, V}[\mu](z) = \frac{\int \exp(\langle K z', Q z \rangle) V z' d\mu(z')}{\int \exp(\langle K z', Q z \rangle) d\mu(z')}$$

→ corresponds to classic definition when  $\mu = \frac{1}{L} \underbrace{\sum_{l=1}^L \sqrt{z_l}}_{i: \hat{\mu}_L}$

If  $z_i \overset{i.i.d.}{\sim} \mu$ , then  $T^{K,Q,V}[\hat{\mu}_L](z) \xrightarrow[L \rightarrow \infty]{a.s.} T^{K,Q,V}[\mu](z)$

infinite prompt regime is often easier to analyze.

Example. Gaussian inputs.

**Lemma** (Castin et al, 2025)

$$\text{If } \mu = \mathcal{N}(m, \Gamma) \text{ then } T^{K,Q,V}[\mu](z) = Vm + V\Gamma K^T Q z$$

Proof

$$T^{K,Q,V}[\mu](z) = V \frac{\int \exp(\langle Kz', Qz \rangle) z' d\mu(z')}{\int \exp(\langle Kz', Qz \rangle) d\mu(z')} = V \mathbb{E}_{z' \sim \nu} [z']$$

with  $\nu$  skewed distribution defined as:  $d\nu(z') = \frac{\exp(\langle Kz', Qz \rangle)}{\int_{z' \sim \mu} \exp(\langle Kz', Qz \rangle)} d\mu(z')$

$\nu$  is a probability distribution and (assuming  $\Gamma$  non-degenerate),

$$\begin{aligned} \frac{d\nu(z')}{dz'} &\propto \exp\left(z'^T K^T Q z - \frac{1}{2} (z' - m)^T \Gamma^{-1} (z' - m)\right) \\ &\propto \exp\left(-\frac{1}{2} (z' - m - \Gamma K^T Q z)^T \Gamma^{-1} (z' - m - \Gamma K^T Q z)\right) \cdot C(z) \end{aligned}$$

$$\Rightarrow \nu = \mathcal{N}(m + \Gamma K^T Q z, \Gamma)$$

$$\text{So } T^{K,Q,V}[\mu](z) = V \mathbb{E}_{z' \sim \nu} [z'] = Vm + V\Gamma K^T Q z \quad \square$$

with Gaussian inputs, softmax attention coincide with (normalized) linear attention in the limit  $L \rightarrow \infty$  where

$$T_{\text{lin}}^{K,Q,V}[\mu](z) = V \int \langle Kz', Qz \rangle z' d\mu(z')$$

→ much easier to analyse (eg optimization dynamics for ICL)

But how do they compare for finite  $L$ ?

## II - Concentration of softmax attention

Prop

If  $\mu$  and  $\nu$  are resp.  $\sigma$  and 1 sub-Gaussian measures with  $\sigma \gg 1$ , then there exists constants  $c_1, c_2 > 0$  depending solely on  $d, V, K, Q$  s.t.

$$\mathbb{E}_{z_1, \dots, z_L \sim \mu^{\otimes L}} \left[ \left\| T_{\text{lin}}^{K,Q,V}[\hat{\mu}_L] - T_{\text{lin}}^{K,Q,V}[\mu] \right\|_{L^2(\nu)}^2 \right] \leq c_1 \frac{\sigma^6 \ln L}{L^{c_2/\sigma^2}}$$

Proof Idea: concentration on both  $\frac{1}{L} \sum_{i=1}^L \exp(\langle Kz_i, Qz \rangle) / V z_i$

and  $\frac{1}{L} \sum_{i=1}^L \exp(\langle Kz_i, Qz \rangle)$

but these are heavy tail random variables.

when  $\max_p \exp(\langle Kz, Qz \rangle) \leq M \rightarrow$  exponential tail concentration

when  $\max_p \exp(\langle Kz, Qz \rangle) > M$  (low proba)  $\rightarrow$  use that  $\|T^{K,Q,V}[\hat{\mu}_L](z)\| \leq \max_p \|Vz\| \square$

Our goal: transfer optimization results derived for linear case ( $L = \infty$ ) to finite  $L$  and softmize attention

For that, we need concentration on the gradients. Use  $U = K^T Q$  parametrization

so that  $T^{U,V} = T^{K,Q,V}$  with

$$T^{U,V}[\hat{\mu}_L](z) = \frac{\int \exp(z^T U z) V z^i d\mu(z)}{\int \exp(z^T U z) d\mu(z)}$$

Under same assumptions + bounded moments of shrunken distribution:

$$\mathbb{E}_{z_1, \dots, z_n \sim \mu^{\otimes n}} \left[ \left\| \nabla_w T^{U,V}[\hat{\mu}_L] - \nabla_w T^{U,V}[\mu_{L(n)}] \right\| \right] \leq c_d \sigma^{\frac{d^2}{2}} \frac{\rho(L)^2}{L^{d/2}}$$

$\rightarrow$  training dynamics should be the same for large values of  $L$ !

### III - In-context linear regression

In-Context Learning (ICL): ability to "learn" task structure directly from the prompt, at inference time

#### Setting (ICL linear regression)

for each new prompt: draw parameter  $w \sim \mathcal{N}(0, I_d)$

and prompt of  $(z_1, \dots, z_n)$  where  $z_i = (x_i, y_i)$  and  $\begin{cases} x_i \sim \mathcal{N}(0, \Sigma) \\ y_i = w^T x_i \end{cases}$

Goal: given a masked query input  $(x_q, 0)$ , predict  $y_q = w^T x_q$ .

Note that given  $w$ ,  $z_p \sim \mathcal{N}(0, \Gamma_w)$  with  $\Gamma_w = \begin{pmatrix} \Sigma & \Sigma w \\ (\Sigma w)^T & \|w\|_\Sigma^2 \end{pmatrix}$

ICL Risk for  $L = \infty$ :

$$\begin{aligned} \mathcal{R}_\infty^{\text{ICL}}(U, V) &= \frac{1}{2} \mathbb{E}_{w \sim \mathcal{N}(0, \Gamma_d)} \left[ \mathbb{E}_{x_q \sim \mathcal{N}(0, \Sigma)} \left[ \left( T_{\mu}^{U, V}(x_q, 0)_{t+1} - w^T x_q \right)^2 \right] \right] \\ &= \frac{1}{2} \mathbb{E}_{w \sim \mathcal{N}(0, \Gamma_d)} \left[ \mathbb{E}_{x_q \sim \mathcal{N}(0, \Sigma)} \left[ \left( U_{t+1}^T \Gamma_w U(x_q, 0)^T - w^T x_q \right)^2 \right] \right] \end{aligned}$$

ICL Risk for  $L < \infty$ :

$$\begin{aligned} \mathcal{R}_L^{\text{ICL}}(U, V) &= \frac{1}{2} \mathbb{E}_w \left[ \mathbb{E}_{x_q} \left[ \left( T_{\hat{\mu}_L}^{U, V}(x_q, 0)_{t+1} - w^T x_q \right)^2 \right] \right] \\ &= \frac{1}{2} \mathbb{E}_w \mathbb{E}_{z_1, \dots, z_L \sim \mathcal{N}(0, \Gamma_i)^{\otimes L}} \left[ \left( \frac{\sum_{t=1}^L \exp(z_t^T U(x_q, 0)^T) \exp(z_t^T z_p)}{\sum_{p=1}^L \exp(z_p^T U(x_q, 0)^T)} - w^T x_q \right)^2 \right] \end{aligned}$$

Then

for a proper init. scheme, then for any  $\epsilon > 0$ ,  $\exists L(\epsilon)$  s.t.  $\forall L \geq L(\epsilon)$ ,

$$\lim_{L \rightarrow \infty} \mathcal{R}_L^{\text{ICL}}(U_L(\mathcal{H}), V_L(\mathcal{H})) \leq \mathcal{R}^{\text{ICL},*} + \epsilon$$

params obtained by gradient flow over  $\mathcal{H}^{\text{ICL}}$

Bayes risk (= instance of label noise)

Idea of proof:

use similarity with the  $L = \infty$  case + the result by Zhang et al (2024)

that provides convergence guarantees to Bayes optimum for linear attention.