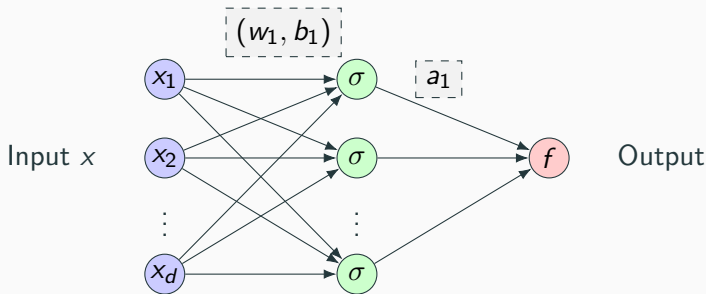


Mean-Field Neural Networks and Transformers

September 17, 2025

Mean-field Neural Networks

Two-Layer Neural Network: Definition

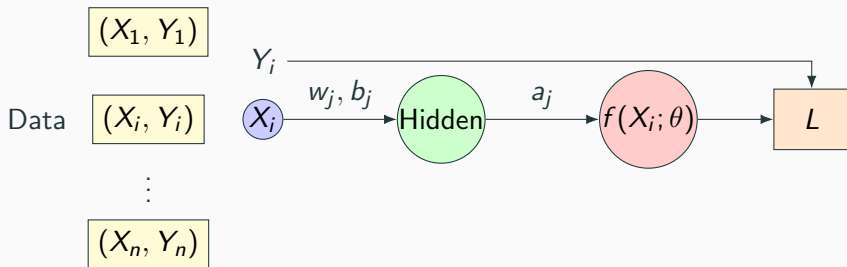


The output $f(x; \theta)$ is given by:

$$f(x; \theta) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\langle w_j, x \rangle + b_j)$$

- $\theta = \{(a_j, w_j, b_j)\}_{j=1}^m$ are the parameters.
- m is the number of neurons.

The Optimization Problem: Non-Convexity



Loss function (e.g., squared loss): $L(\theta) = \sum_{i=1}^n (Y_i - f(X_i; \theta))^2$

$$L(\theta) = \sum_{i=1}^n \left(Y_i - \frac{1}{m} \sum_{j=1}^m a_j \sigma(\langle w_j, X_i \rangle + b_j) \right)^2$$

This loss function $L(\theta)$ is highly **non-convex**.

Lifting to Measure Space



$$\text{Let } \mu = \frac{1}{m} \sum_{j=1}^m \delta_{(a_j, w_j, b_j)}.$$

Network output becomes an integral:

$$f(x; \mu) = \int_{\Omega} a\sigma(\langle w, x \rangle + b) \mu(d\omega) = \int_{\Omega} \rho(x; \omega) \mu(d\omega)$$

The loss becomes a functional on the space of measures:

$$L : \mathcal{P}_2(\Omega) \rightarrow \mathbb{R}$$

$$L(\mu) = \sum_{i=1}^n \left(Y_i - \int_{\Omega} \rho(X_i; \omega) \mu(d\omega) \right)^2$$

Gradient Flows: Euclidean vs. Wasserstein

Correspondence Proposition

The Wasserstein gradient flow $(\mu_t)_{t \geq 0}$ of $L(\mu)$, when initialised at an empirical measure $\mu_0 = \mu_\theta$ (where $\theta = \{(a_j, w_j, b_j)\}_{j=1}^m$ and $\mu_\theta = \frac{1}{m} \sum_{j=1}^m \delta_{(a_j, w_j, b_j)}$), satisfies $\mu_t = \mu_{\theta_t}$ for all $t \geq 0$. Here $(\theta_t)_{t \geq 0}$ is the (time-rescaled) Euclidean gradient flow of $L(\theta)$ initialised at θ .

| | |
|--|--|
| State Represents: θ_t | State Represents: $\mu_t = \mu_{\theta_t}$ |
| Loss Function: $L(\theta_t)$ | Loss Functional: $L(\mu_t) = L(\mu_{\theta_t})$ |
| Optimization Dynamics: Euclidean Gradient Flow $\dot{\theta}_t = -\nabla_\theta L(\theta_t)$ | Optimization Dynamics: Wasserstein Gradient Flow (PDE) $\partial_t \mu_t = \nabla \cdot \left(\mu_t \nabla \frac{\delta L}{\delta \mu_t} \right)$ |

Geodesic Non-Convexity Remains!

The lifted loss functional $L(\mu)$:

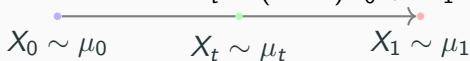
$$L(\mu) = \sum_{i=1}^n \left(Y_i - \int_{\Omega} \rho(x; \omega) \mu(d\omega) \right)^2$$

- $L(\mu)$ is "convex" for linear combinations $\lambda\mu_1 + (1 - \lambda)\mu_2$.
- But convex combinations (geodesics) in Wasserstein space $\mathcal{P}_2(\Omega)$ are defined via optimal transport, not linear mixing of measures.
- $L(\mu)$ is generally **not** geodesically convex.

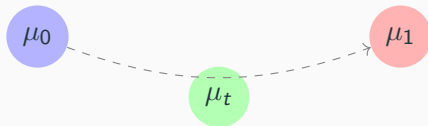
Geodesic Non-Convexity Remains!

Wasserstein Geodesics vs. Linear Interpolation:

Particle Path $X_t = (1 - t)X_0 + tX_1$



The diagram shows a horizontal line with an arrow pointing from left to right. Three points are marked on the line: a blue dot on the left, a green dot in the middle, and a red dot on the right. Below each dot is a label: $X_0 \sim \mu_0$ under the blue dot, $X_t \sim \mu_t$ under the green dot, and $X_1 \sim \mu_1$ under the red dot.



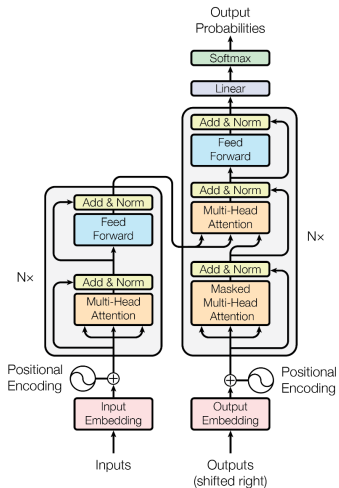
Advantages of the Measure Perspective

Allows considering measures other than the empirical $\frac{1}{m} \sum_{i=1}^m \delta_{\omega_j}$ and provides analytical benefits:

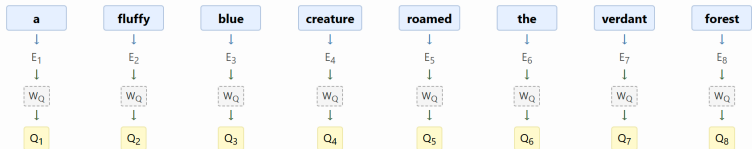
- **Chizat & Bach '18, "On the Global Convergence of Gradient Descent for Over-parameterized Models using Optimal Transport"**
- **Mei et al. '18, "A Mean Field View of the Landscape of Two-Layers Neural Networks"**

Transformers

Transformer Architecture



Self-Attention Example: QKV Calculation



| | Q ₁ | Q ₂ | Q ₃ | Q ₄ | Q ₅ | Q ₆ | Q ₇ | Q ₈ |
|---|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|
| a → E ₁ → w _k → K ₁ | K ₁ ·Q ₁ | K ₁ ·Q ₂ | K ₁ ·Q ₃ | K ₁ ·Q ₄ | K ₁ ·Q ₅ | K ₁ ·Q ₆ | K ₁ ·Q ₇ | K ₁ ·Q ₈ |
| fluffy → E ₂ → w _k → K ₂ | K ₂ ·Q ₁ | K ₂ ·Q ₂ | K ₂ ·Q ₃ | K ₂ ·Q ₄ | K ₂ ·Q ₅ | K ₂ ·Q ₆ | K ₂ ·Q ₇ | K ₂ ·Q ₈ |
| blue → E ₃ → w _k → K ₃ | K ₃ ·Q ₁ | K ₃ ·Q ₂ | K ₃ ·Q ₃ | K ₃ ·Q ₄ | K ₃ ·Q ₅ | K ₃ ·Q ₆ | K ₃ ·Q ₇ | K ₃ ·Q ₈ |
| creature → E ₄ → w _k → K ₄ | K ₄ ·Q ₁ | K ₄ ·Q ₂ | K ₄ ·Q ₃ | K ₄ ·Q ₄ | K ₄ ·Q ₅ | K ₄ ·Q ₆ | K ₄ ·Q ₇ | K ₄ ·Q ₈ |
| roamed → E ₅ → w _k → K ₅ | K ₅ ·Q ₁ | K ₅ ·Q ₂ | K ₅ ·Q ₃ | K ₅ ·Q ₄ | K ₅ ·Q ₅ | K ₅ ·Q ₆ | K ₅ ·Q ₇ | K ₅ ·Q ₈ |
| the → E ₆ → w _k → K ₆ | K ₆ ·Q ₁ | K ₆ ·Q ₂ | K ₆ ·Q ₃ | K ₆ ·Q ₄ | K ₆ ·Q ₅ | K ₆ ·Q ₆ | K ₆ ·Q ₇ | K ₆ ·Q ₈ |
| verdant → E ₇ → w _k → K ₇ | K ₇ ·Q ₁ | K ₇ ·Q ₂ | K ₇ ·Q ₃ | K ₇ ·Q ₄ | K ₇ ·Q ₅ | K ₇ ·Q ₆ | K ₇ ·Q ₇ | K ₇ ·Q ₈ |
| forest → E ₈ → w _k → K ₈ | K ₈ ·Q ₁ | K ₈ ·Q ₂ | K ₈ ·Q ₃ | K ₈ ·Q ₄ | K ₈ ·Q ₅ | K ₈ ·Q ₆ | K ₈ ·Q ₇ | K ₈ ·Q ₈ |

Self-Attention Example: Scores and Normalization

Unnormalized
Attention Pattern

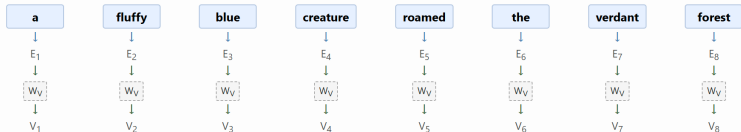
| | | | | | |
|-------|-------|-------|-------|-------|-------|
| +3.53 | +0.80 | +1.96 | +4.48 | +3.74 | -1.95 |
| -∞ | -0.30 | -0.21 | +0.82 | +0.29 | +2.91 |
| -∞ | +0.89 | +0.67 | +2.99 | -0.41 | |
| -∞ | -∞ | +1.31 | +1.73 | -1.48 | |
| -∞ | -∞ | -∞ | +3.07 | +2.94 | |
| -∞ | -∞ | -∞ | -∞ | +0.31 | |

softmax

Normalized
Attention Pattern

| | | | | | |
|------|------|------|------|------|------|
| 1.00 | 0.75 | 0.69 | 0.92 | 0.46 | 0.00 |
| 0.00 | 0.25 | 0.08 | 0.02 | 0.01 | 0.46 |
| 0.00 | 0.00 | 0.24 | 0.02 | 0.22 | 0.02 |
| 0.00 | 0.00 | 0.00 | 0.04 | 0.06 | 0.01 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.24 | 0.48 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.03 |

Self-Attention Example: Weighted Value Sum



| $a \rightarrow E_1 \rightarrow W_v \rightarrow V_1$ | 1.00 V_1 | 0.00 V_1 | 0.00 V_1 | 0.00 V_1 | 0.00 V_1 | 0.00 V_1 | 0.00 V_1 | 0.00 V_1 |
|--|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $fluffy \rightarrow E_2 \rightarrow W_v \rightarrow V_2$ | 0.00 V_2 | 1.00 V_2 | 0.00 V_2 | 0.42 V_2 | 0.00 V_2 | 0.00 V_2 | 0.00 V_2 | 0.00 V_2 |
| $blue \rightarrow E_3 \rightarrow W_v \rightarrow V_3$ | 0.00 V_3 | 0.00 V_3 | 1.00 V_3 | 0.58 V_3 | 0.00 V_3 | 0.00 V_3 | 0.00 V_3 | 0.00 V_3 |
| $creature \rightarrow E_4 \rightarrow W_v \rightarrow V_4$ | 0.00 V_4 | 0.00 V_4 | 0.00 V_4 | 0.00 V_4 | 0.00 V_4 | 0.00 V_4 | 0.00 V_4 | 0.00 V_4 |
| $roamed \rightarrow E_5 \rightarrow W_v \rightarrow V_5$ | 0.00 V_5 | 0.00 V_5 | 0.00 V_5 | 0.00 V_5 | 0.00 V_5 | 0.01 V_5 | 0.00 V_5 | 0.00 V_5 |
| $the \rightarrow E_6 \rightarrow W_v \rightarrow V_6$ | 0.00 V_6 | 0.00 V_6 | 0.00 V_6 | 0.00 V_6 | 0.99 V_6 | 1.00 V_6 | 0.00 V_6 | 0.00 V_6 |
| $verdant \rightarrow E_7 \rightarrow W_v \rightarrow V_7$ | 0.00 V_7 | 0.00 V_7 | 0.00 V_7 | 0.00 V_7 | 0.00 V_7 | 0.00 V_7 | 1.00 V_7 | 0.00 V_7 |
| $forest \rightarrow E_8 \rightarrow W_v \rightarrow V_8$ | 0.00 V_8 | 0.00 V_8 | 0.00 V_8 | 0.00 V_8 | 0.00 V_8 | 0.00 V_8 | 0.00 V_8 | 1.00 V_8 |
| | Σ | Σ | Σ | Σ | Σ | Σ | Σ | Σ |
| | \downarrow | \downarrow | \downarrow | \downarrow | \downarrow | \downarrow | \downarrow | \downarrow |
| | ΔE_1 | ΔE_2 | ΔE_3 | ΔE_4 | ΔE_5 | ΔE_6 | ΔE_7 | ΔE_8 |

Self-Attention as Dynamics (1/4): Iterative Scheme

Consider applying the attention update iteratively, like residual connections in deep networks:

$$x_{t+1}^i = x_t^i + \left(V \frac{\sum_{j=1}^N x_t^j e^{\langle Qx_t^i, Kx_t^j \rangle}}{\sum_{l=1}^N e^{\langle Qx_t^i, Kx_t^l \rangle}} \right)$$

Self-Attention as Dynamics (2/4): Continuous ODE

Consider applying the attention update iteratively, like residual connections in deep networks:

$$\dot{x}_t^i = V \frac{\sum_{j=1}^N x_t^j e^{\langle Qx_t^i, Kx_t^j \rangle}}{\sum_{l=1}^N e^{\langle Qx_t^i, Kx_t^l \rangle}}$$

Self-Attention as Dynamics (3/4): Mean-Field Limit

Consider applying the attention update iteratively, like residual connections in deep networks:

$$\dot{x}_t^i = V \frac{\int y e^{\langle Qx_t^i, Ky \rangle} \mu_t(dy)}{\int e^{\langle Qx_t^i, Ky \rangle} \mu_t(dy)}$$

Assume $Q = K = V = I$ for simplicity:

$$\dot{x}_t^i = \frac{\int y e^{\langle x_t^i, y \rangle} \mu_t(dy)}{\int e^{\langle x_t^i, y \rangle} \mu_t(dy)} = \nabla_x \left[\log \int e^{\langle x, y \rangle} \mu_t(dy) \right]_{x=x_t^i}$$

Self-Attention as Dynamics (4/4): WGF Connection?

Recall: Particle Interpretation of WGF:

$$\dot{X}_t = -\nabla \frac{\delta \mathcal{F}}{\delta \mu}(X_t)$$

From the previous slide (with $Q = K = V = I$):

$$\dot{x}_t^i = \nabla \underbrace{\left[\log \int e^{\langle x, y \rangle} \mu_t(dy) \right]}_{\Psi(x; \mu_t)}(x_t^i)$$

Question

$\Psi(x; \mu_t) \stackrel{?}{=} -\frac{\delta \mathcal{F}}{\delta \mu_t}(x)$ for some energy functional \mathcal{F} describing the self-attention dynamics?

No!

- **Reference:** Sander et al. '22, "Sinkformers: Transformers with Doubly Stochastic Attention".
- **Issue:** Asymmetry.

Unnormalised Self-Attention (1/2)

Remedy: What if we remove the denominator (the Softmax normalization)?

Consider the unnormalised dynamics (assuming $Q = K = V = I$ for now):

$$\dot{x}_t^i = \int y e^{\langle x_t^i, y \rangle} \mu_t(dy) = \nabla_x \left[\int e^{\langle x, y \rangle} \mu_t(dy) \right]_{x=x_t^i}$$

Unnormalised Self-Attention (2/2)

Remedy: What if we remove the denominator (the Softmax normalization)?

Consider the unnormalised dynamics (assuming $Q = K = V = I$ for now):

$$\dot{x}_t^i = \int y e^{\langle x_t^i, y \rangle} \mu_t(dy) = \nabla_x \left[\underbrace{\int e^{\langle x, y \rangle} \mu_t(dy)}_{\delta \mathcal{F}(\mu_t)} \right]_{x=x_t^i}$$

where

$$\mathcal{F}(\mu) = - \int \int e^{\langle x, y \rangle} \mu(dx) \mu(dy)$$

Unnormalised Attention as WGF

Proposition

The unnormalised self-attention dynamics (with $Q = K = V = I$, after time rescaling):

$$\dot{x}_t^i = \int y e^{\langle x_t^i, y \rangle} \mu_t(dy)$$

is the Wasserstein gradient flow of the interaction energy:

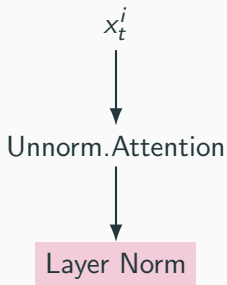
$$\mathcal{F}(\mu) = - \int \int e^{\langle x, y \rangle} \mu(dx) \mu(dy)$$

defined on $\mathcal{P}_2(\mathbb{R}^d)$.

Implication: Divergence

$\mathcal{F}(\delta_z) = -e^{\langle z, z \rangle} = -e^{\|z\|^2}$. As $\|z\| \rightarrow \infty$, $\mathcal{F}(\delta_z) \rightarrow -\infty$.

Adding Layer Normalisation



Corresponds to the dynamics

$$\dot{x}_t^i = P_{x_t^i} \left(\int y e^{\langle x_t^i, y \rangle} \mu_t(dy) \right)$$

where P_x is the projection onto $T_x(S^{d-1})$

Layer-Normalised Unnormalised Attention as WGF

Proposition

The layer-normalised unnormalised self-attention dynamics:

$$\dot{x}_t^i = P_{x_t^i} \left(\int y e^{\langle x_t^i, y \rangle} \mu_t(dy) \right)$$

is the Wasserstein gradient flow of the same interaction energy:

$$\mathcal{F}(\mu) = - \int \int e^{\langle x, y \rangle} \mu(dx) \mu(dy)$$

but now considered on the space of probability measures on the sphere, $\mathcal{P}_2(S^{d-1})$.

Implication: Concentration

Minimizers are Dirac masses δ_z for $z \in S^{d-1}$

The dynamics tend to form a single cluster on the sphere.

The interaction energy on the sphere:

$$\mathcal{F}(\mu) = - \int_{S^{d-1}} \int_{S^{d-1}} e^{\langle x, y \rangle} \mu(dx) \mu(dy)$$

- This functional $\mathcal{F}(\mu)$ is **not** geodesically convex on $\mathcal{P}_2(S^{d-1})$.
- It admits many stationary points where the Wasserstein gradient vanishes

The interaction energy on the sphere:

$$\mathcal{F}(\mu) = - \int_{S^{d-1}} \int_{S^{d-1}} e^{\langle x, y \rangle} \mu(dx) \mu(dy)$$

- **Geshkovski et al. '24, "A Mathematical Perspective on Transformers"**: These points are in fact saddle points, guaranteeing asymptotic convergence to a single cluster when dynamics are initialized in a generic position.