

A Theoretical Framework For Overfitting In Energy-based Modeling

Giovanni Catania, Aurélien Decelle, Cyril Furtlehner, Beatriz Seoane

Universidad Complutense Madrid (ES)

International Conference on Machine Learning, 2025

Motivation

Study the impact of the amount data in the training of Energy-Based Models (EBMs) and how overfitting emerges when data is limited

In generative models, overfitting occurs when the model “memorizes” the training data instead of learning the underlying data distribution.

- poor diversity in generated samples, lack of variability present in real data
- model learns specific noise-dominated information in the data

Energy based models (EBMs) in generative AI

EBMs encode the empirical distribution of a dataset into a Boltzmann distribution with a given Energy function



Data

Model

$$p_{\text{data}}(\mathbf{x}) \sim p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{Z_{\boldsymbol{\theta}}} e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}$$

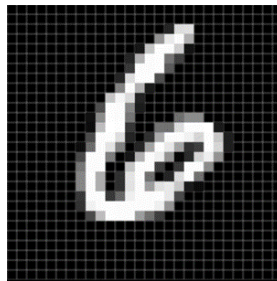
Boltzmann distribution



rooted in statistical physics

Dataset (e.g. MNIST)

Datum \mathbf{x}



$\boldsymbol{\theta}$: vector of parameters to be trained

Training: log-likelihood maximization

Used for generative purposes and for interpretability of the effective model

Theoretical analysis of overfitting in EBM

- Use a simple model (analytically solvable) for a synthetic experiment:
- Track the quality of the inferred model as a function of the number of samples

↓
Gaussian Model

Gaussian Energy-based Model (GEBM)

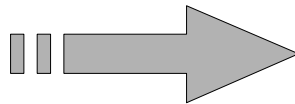
True model

$$E_{\theta}(\mathbf{x}) = \frac{1}{2} \sum_{ij} J_{ij} x_i x_j$$

$\mathbf{x} \in \mathbb{R}$

$$\mathbf{J} = (\mathbf{C}^*)^{-1}$$

Population
covariance matrix



Sample M configurations from

$$^*p_{\mathbf{J}}(\mathbf{x}) = \frac{1}{Z_{\mathbf{J}}} e^{-E_{\mathbf{J}}(\mathbf{x})}$$

$$\hat{\mathbf{C}}^M = \frac{1}{M} \sum_{\mu=1}^M \mathbf{x}_{\mu} \mathbf{x}_{\mu}^T$$

Empirical
covariance matrix

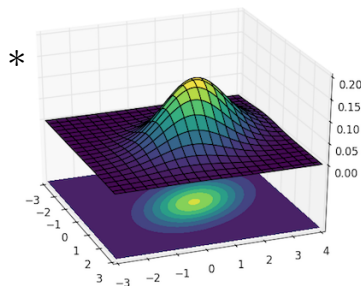


Infer back the model
from these M samples

$$\hat{\mathbf{C}}^M \xrightarrow{M \rightarrow \infty} \mathbf{C}^*$$

$M \equiv$ Number of data

Multivariate Gaussian



Gaussian Energy-based Model (GEBM)

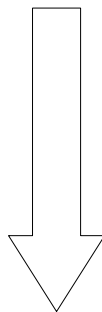
Why Gaussian Model?

Analytically solvable, both Maximum likelihood estimator and (most importantly) the training dynamics

Max-likelihood (ML) estimator

$$\mathbf{J}^{\text{ML}} = (\hat{\mathbf{C}}^M)^{-1}$$

same eigenvector basis



$$\mathbf{J}(t) = \sum_{\alpha} \mathbf{J}_{\alpha}(t) \mathbf{u}_{\alpha} \mathbf{u}_{\alpha}^{\top}$$

Study training dynamics (likelihood maximization)
by projecting on eigenvector basis
→ each eigenvalue of \mathbf{J} now evolves independently on the others

Separation of learning timescales

Eigenvalue of $\frac{dJ_\alpha}{dt}$ Eigenvalue of C^M

$$\frac{dJ_\alpha}{dt} = -c_\alpha^M + \frac{1}{J_\alpha}$$

Separation of time-scales

Modes corresponding to stronger correlations are learnt faster

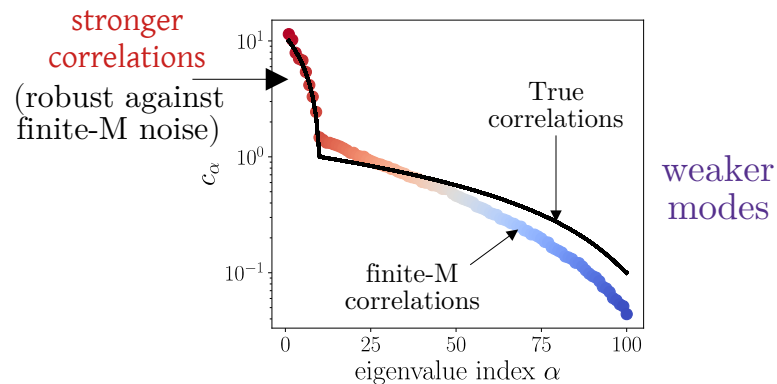
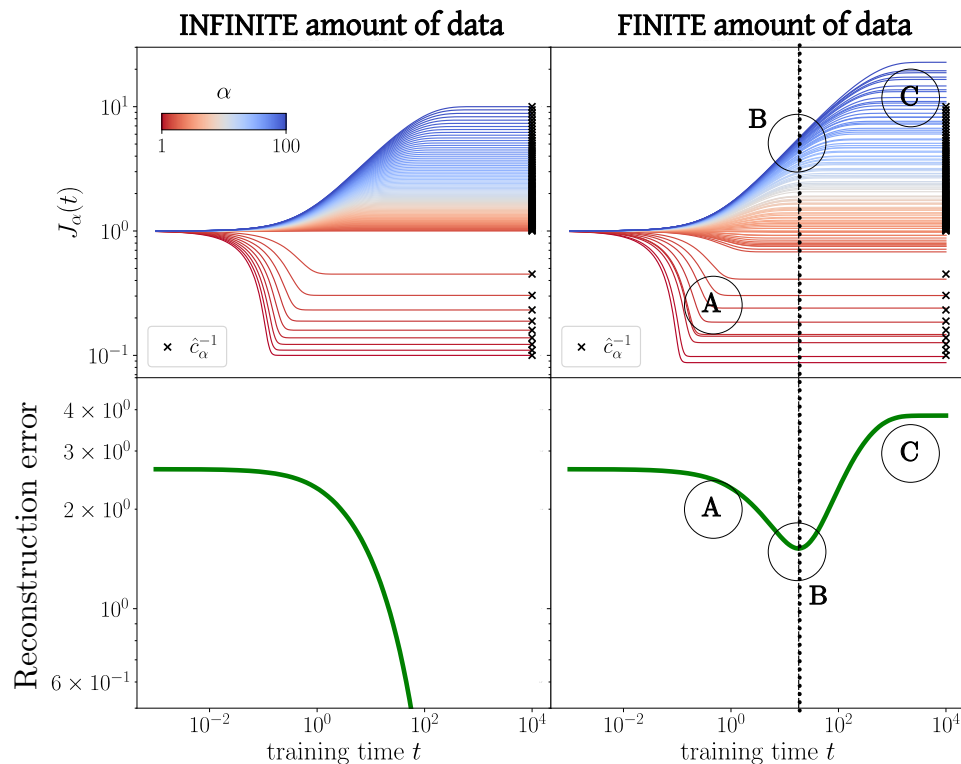
$$J_\alpha(t) = \frac{1}{c_\alpha^M} + \frac{1}{c_\alpha^M} W \left(\text{const } e^{-(c_\alpha^M)^2 t} \right)$$

1. Model learns information sequentially way*
stronger PCA direction first, then weaker ones

$$\text{Learning time-scale} \propto \frac{1}{(c_\alpha^M)^2}$$

2. Strong/weak PCA modes have very different fluctuations w.r.t. the number of data

Eigenvalues' evolution



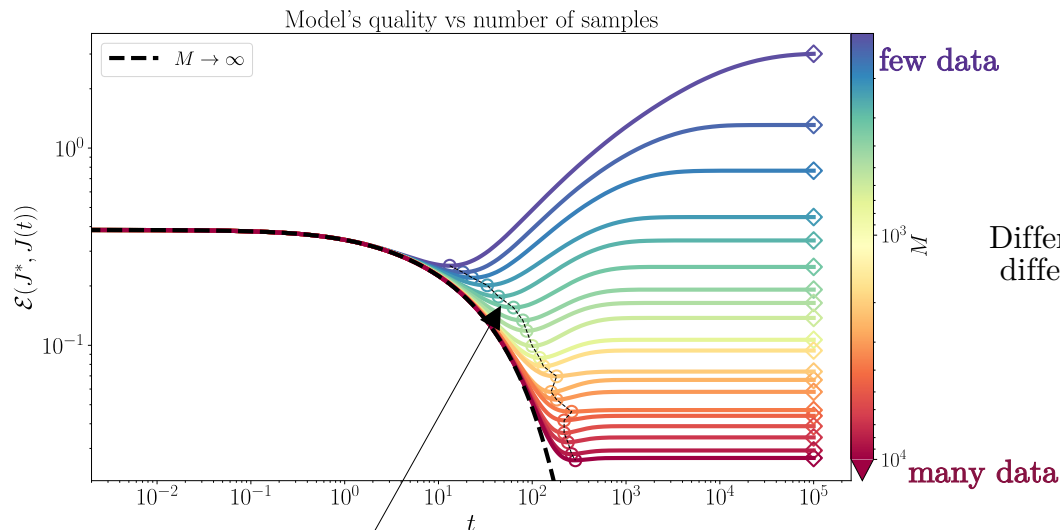
- A) Eigenvalues of **stronger modes** are the first to converge
→ error decreases
- B) **weaker modes** are starting to be learnt
→ error decreases (up to minimum)
eigenvalues are closer to the true value than to the fixed point!
- C) **weaker modes** converge to the fixed point,
error increases after minimum is reached
→ error dominated by fluctuations of weaker correlations due to the low number of samples

$$\mathcal{E}_J = \|\mathbf{J}^{\text{true}} - \mathbf{J}(t)\|$$

Early stopping points in training dynamics

Non-monotonic behavior (w.r.t. training time) of discrepancy between true and inferred model

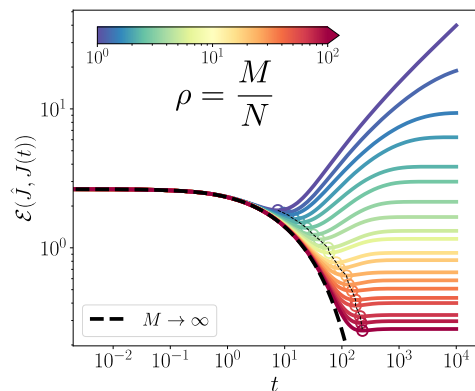
Models inferred with few training data are **worse** at fixed point than during training



early-stopping point at which best inference is achieved, after which model starts to overfit

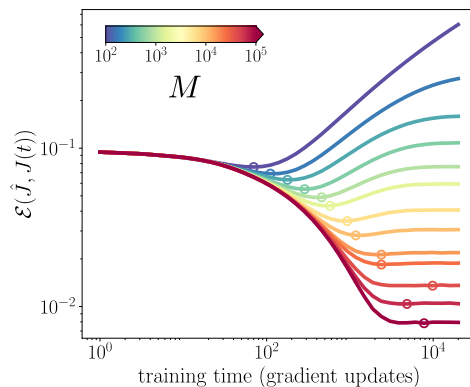
Different EBM, same phenomenology

Gaussian Model
(GEBM)



Boltzmann Machine
Inverse Ising

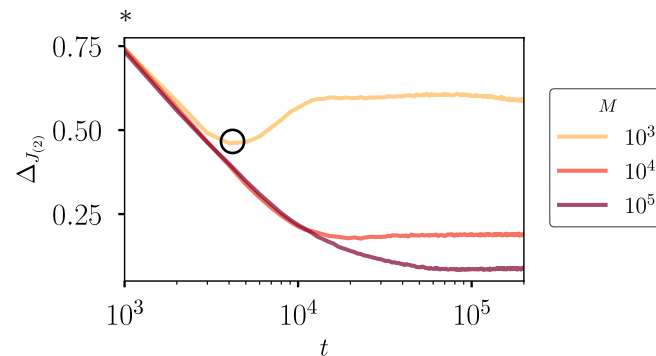
Dataset: equilibrium configurations
from 2D Ising model (high-Temp)



Similar analysis of training dynamics
can be carried out analytically,
using Mean-Field approximation

Restricted Boltzmann Machine

Dataset: equilibrium configurations
from 1D Ising model at high T



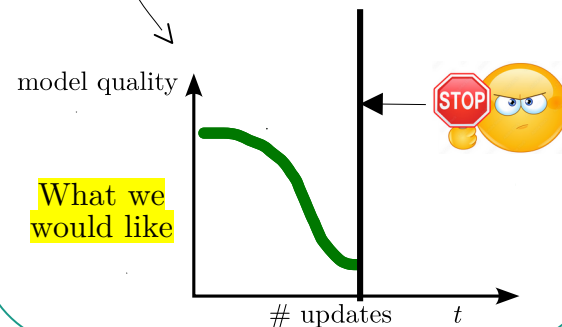
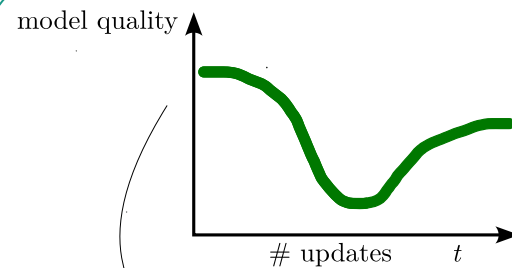
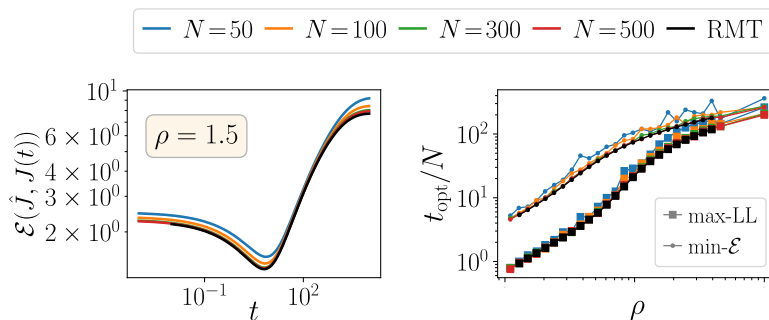
*Taken from
Decelle, Furtlehner, Navas Gómez, Seoane,
SciPost Physics 16(4)095 (2024)

Random matrix theory analysis

Asymptotic analysis through Random Matrix theory (RMT) to analyze finite-samples fluctuations in the training dynamics

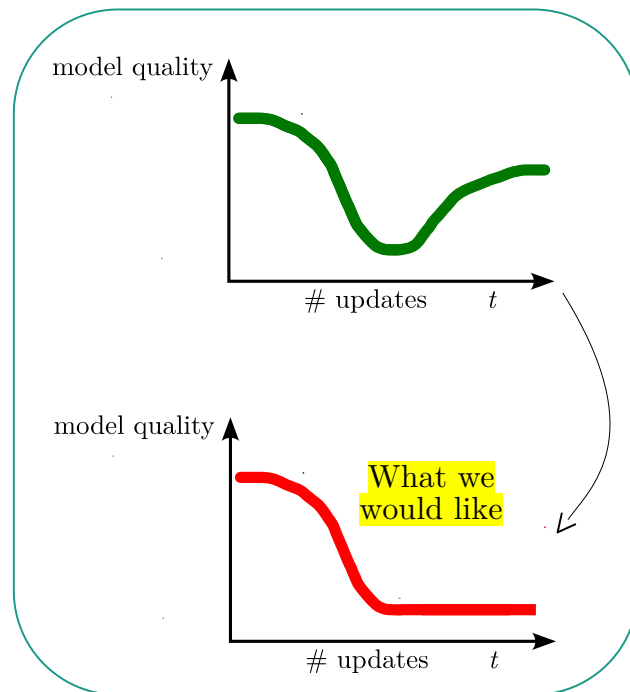
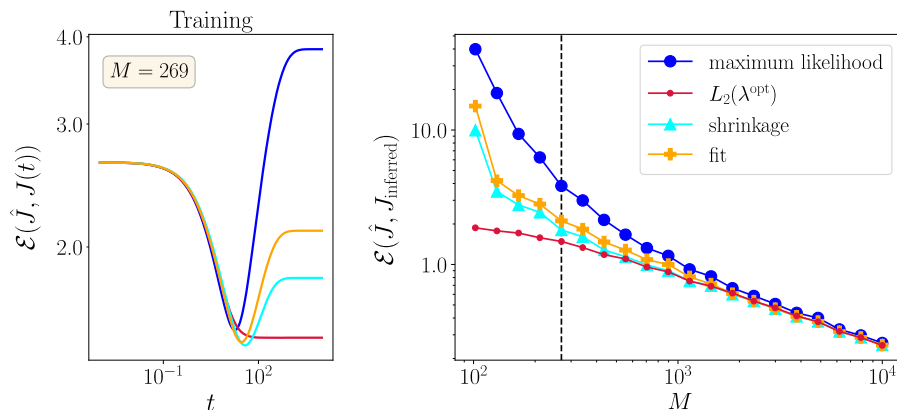
Exact on GEBM

$$\begin{aligned} N &\rightarrow \infty \\ M &\rightarrow \infty \\ \rho &= \frac{M}{N} \text{ finite} \end{aligned}$$



Protocols to mitigate overfitting

- regularization priors
- shrinkage correction protocols
- downsampling-based modes fitting



Extensions to more complex EBM

Study of overfitting in arbitrary complex EBMs can be done using the score-matching algorithm

$$p_{\theta}(\mathbf{x}) = \frac{1}{Z_{\theta}} e^{-E_{\theta}(\mathbf{x})}$$

Learning dynamics of score function governed by a Neural Tangent Kernel

$$\frac{d\psi(\mathbf{x}|\theta_t)}{dt} = -\hat{\mathbb{E}}_{x'} \left[K_t(\mathbf{x}, \mathbf{x}') \psi(\mathbf{x}'|\theta_t) \right] + \hat{\phi}_t(\mathbf{x})$$

Score function

$$\psi(\mathbf{x}|\theta) = -\nabla_{\mathbf{x}} E(\mathbf{x}|\theta)$$

Similar learning dynamics to GEBM w.r.t. empirical covariance of latent feature in the tangent space \rightarrow can lead to similar mechanism that justify the onset of overfitting

Summary

- Introduction of a novel theoretical framework to study overfitting in EBM
- **Interplay between learning timescales** associated to different PCA directions (with different finite-sample fluctuations) can result in overfitting.
- Analysis on GEBM, asymptotics through RMT
- Theoretical extension on Boltzmann Machine (high-T), **extension to generic EBM in the context of NTK of the score function dynamics**
- sets the stage for
 - a) early-stopping point determination through RMT
 - b) extension of data-correction protocols to non-pairwise EBMs

arXiv: 2501.19158

ICML POSTER ID: 45237



with



Aurélien Decelle
Universidad Politecnica Madrid (ES)



Cyril Furtlehner
INRIA, Université Paris Saclay (FR)



Beatriz Seoane
Universidad Complutense Madrid (ES)

Thank you