

Maria Han Veiga

Ohio Stae University, USA

Mathematical aspects of generative machine learning models

Abstract:

Diffusion probabilistic models have become the state-of-the-art tool in generative methods, used to generate high-resolution samples from very high-dimension distributions (e.g. images). It relies on a forward-time stochastic differential equation (SDE) that slowly injects noise on the data distribution, transforming it into a known prior distribution, and a reverse-time SDE that transforms the prior distribution back into the data distribution by removing the noise. The reverse-time SDE follows the time-dependent gradient field (score) of the perturbed data distribution.

Although very effective, there are some drawback to this method: 1) as opposed to variational encoders, the dimension of the problem remains high during the generation process and 2) they can be prone to memorisation of the dataset. In this talk, we consider the second point: the learned score can overfit the finite dataset, making the reverse-time SDE sample mostly training points. In our work, we interpret the learned empirical score as a noisy approximation of the true score and study its covariance matrix, showing that in high-dimension for small time, the noise variance blows up. To reduce this variance, we introduce a smoothing operator on the empirical score and analyse its bias-variance properties. Consequences of this approach are better generalisation of the dataset and thus, a reduction on the amount of samples required to approximate a distribution.

room 40.03.003 (Emil Fischer Str. 40)

Thursday, May. 22, 2025 at 12:30 pm

Zu diesem Vortrag sind Sie herzlich eingeladen.