

# Drifting Models

Yutao Chen

✘ May 19 2026

✘ May 19 2026

## Contents

Drifting Field .....	1
Theoretical Connections .....	2
Score Matching .....	2
Reverse KL Divergence .....	2
Kernel-Gradient Drifting .....	3

Drifting models (Deng et al., 2026) are generative models that admits fast one-step sampling by minimizing *drifting* at train time.

## Drifting Field

Let  $p_{\text{data}}(\mathbf{x})$  be the data distribution, and  $q_{\theta}(\mathbf{x})$  be the model distribution induced by a neural network  $\mathbf{x} = f_{\theta}(\mathbf{z})$  with noise  $\mathbf{z} \sim p(\mathbf{z})$ . Drifting models approximate  $p(\mathbf{x})$  with  $q_{\theta}(\mathbf{x})$  by minimizing

$$\mathcal{L}(\theta) = \mathbb{E}_{\mathbf{x} \sim q_{\theta}} [\|\mathbf{x} - \text{sg}(\mathbf{x} + V_{p,q}(\mathbf{x}))\|^2],$$

where  $\text{sg}(\cdot)$  is the stop-gradient operator, and  $V_{p,q}(\mathbf{x})$  is the drifting field.

The drifting field  $V_{p,q}(\mathbf{x})$  for distributions  $p, q$  on  $\mathbb{R}^d$  is defined as

$$V_{p,q}(\mathbf{x}) = V_p(\mathbf{x}) - V_q(\mathbf{x}),$$

with

$$V_p(\mathbf{x}) = \frac{\mathbb{E}_{\mathbf{y} \sim p}[k(\mathbf{x}, \mathbf{y})\mathbf{y}]}{\mathbb{E}_{\mathbf{y} \sim p}[k(\mathbf{x}, \mathbf{y})]} - \mathbf{x}, \quad V_q(\mathbf{x}) = \frac{\mathbb{E}_{\mathbf{y} \sim q}[k(\mathbf{x}, \mathbf{y})\mathbf{y}]}{\mathbb{E}_{\mathbf{y} \sim q}[k(\mathbf{x}, \mathbf{y})]} - \mathbf{x},$$

where  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a kernel function. Intuitively,  $V_p(\mathbf{x}) / V_q(\mathbf{x})$  can be interpreted as *attraction / repulsion* forces respectively, where the former pulls  $q_{\theta}(\mathbf{x})$  towards  $p(\mathbf{x})$  and the latter disperses  $q_{\theta}$ .

In practice, we approximate  $V_{p,q}(\mathbf{x})$  empirically by drawing batches of attracting samples  $\{\mathbf{y}^+\} \sim p$  and repulsing samples  $\{\mathbf{y}^-\} \sim q$ .

## Theoretical Connections

### Score Matching

The drifting field  $V_{p,q}(\mathbf{x})$  can be interpreted as an approximation of the score mismatch for densities  $p, q$  with *Gaussian* kernels.

Specifically, consider Gaussian kernels with temperature  $\tau > 0$

$$k(\mathbf{x}, \mathbf{y}) = \frac{1}{(\sqrt{2\pi}\tau)^d} \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{2\tau^2}\right).$$

The kernel smoothed (convolved) density of  $p(\mathbf{x})$  on  $\mathbb{R}^d$  is defined as

$$\tilde{p}(\mathbf{x}) := (p * k)(\mathbf{x}) = \mathbb{E}_{\mathbf{y} \sim p}[k(\mathbf{x}, \mathbf{y})].$$

For Gaussian kernels  $k$ , we can show that  $V_{p,q}(\mathbf{x})$  is proportional to the score mismatch between the kernel smoothed densities  $\tilde{p}, \tilde{q}$

$$V_{p,q}(\mathbf{x}) \propto \nabla_{\mathbf{x}} \log \tilde{p}(\mathbf{x}) - \nabla_{\mathbf{x}} \log \tilde{q}(\mathbf{x}).$$

*Proof.* Note that  $\nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{y}) = \frac{1}{\tau^2} k(\mathbf{x}, \mathbf{y})(\mathbf{y} - \mathbf{x})$ , and that

$$\nabla_{\mathbf{x}} \log \tilde{p}(\mathbf{x}) = \frac{\nabla_{\mathbf{x}} \tilde{p}(\mathbf{x})}{\tilde{p}(\mathbf{x})}, \quad \nabla_{\mathbf{x}} \tilde{p}(\mathbf{x}) = \mathbb{E}_{\mathbf{y} \sim p}[\nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{y})].$$

Therefore, we have

$$\nabla_{\mathbf{x}} \log \tilde{p}(\mathbf{x}) = \frac{1}{\tau^2} \frac{\mathbb{E}_{\mathbf{y} \sim p}[k(\mathbf{x}, \mathbf{y})(\mathbf{y} - \mathbf{x})]}{\mathbb{E}_{\mathbf{y} \sim p}[k(\mathbf{x}, \mathbf{y})]} = \frac{1}{\tau^2} V_p(\mathbf{x}).$$

Similarly, we can show that  $\nabla_{\mathbf{x}} \log \tilde{q}(\mathbf{x}) = \frac{1}{\tau^2} V_q(\mathbf{x})$ . Therefore, we have

$$V_{p,q}(\mathbf{x}) = V_p(\mathbf{x}) - V_q(\mathbf{x}) = \tau^2 (\nabla_{\mathbf{x}} \log \tilde{p}(\mathbf{x}) - \nabla_{\mathbf{x}} \log \tilde{q}(\mathbf{x})). \quad \blacksquare$$

### Reverse KL Divergence

Consider the that the Wasserstein gradient flow of the reverse KL divergence  $\mathcal{F}(q) = \mathbb{D}_{\text{KL}}(q \parallel p)$ . The functional derivative w.r.t.  $q(\mathbf{x})$  is

$$\frac{\partial \mathcal{F}}{\partial q}(\mathbf{x}) = 1 + \log q(\mathbf{x}) - \log p(\mathbf{x}),$$

and the corresponding gradient field w.r.t.  $\mathbf{x}$  is given by

$$\nabla_{\mathbf{x}} \frac{\partial \mathcal{F}}{\partial q}(\mathbf{x}) = \nabla_{\mathbf{x}} \log q(\mathbf{x}) - \nabla_{\mathbf{x}} \log p(\mathbf{x}).$$

We can see the drifting field  $V_{p,q}(\mathbf{x})$  yields the steepest *descent* direction of  $\mathbb{D}_{\text{KL}}(\tilde{q} \parallel \tilde{p})$  for kernel smoothed densities  $\tilde{p}, \tilde{q}$ .

When assuming kernels are *characteristic*, minimizing  $\mathbb{D}_{\text{KL}}(\tilde{q} \parallel \tilde{p})$  is equivalent to minimizing  $\mathbb{D}_{\text{KL}}(q \parallel p)$ .

We refer readers to [Gretton et al. \(2026\)](#) for further discussions.

### Kernel-Gradient Drifting

Note that the connection between the drifting field and score mismatch

$$V_{p,q}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x}) - \nabla_{\mathbf{x}} \log q(\mathbf{x})$$

only holds for Gaussian kernels, which satisfy  $\nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{y}) \propto (\mathbf{y} - \mathbf{x})$ .

For general kernels, we consider replacing  $(\mathbf{y} - \mathbf{x})$  with  $\nabla_{\mathbf{x}} \log k(\mathbf{x}, \mathbf{y})$ , yielding the kernel-gradient drifting in [Esteban-Casadevall et al. \(2026\)](#):

$$V_{p,q}(\mathbf{x}) = \frac{\mathbb{E}_{\mathbf{y} \sim p}[\nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{y})]}{\mathbb{E}_{\mathbf{y} \sim p}[k(\mathbf{x}, \mathbf{y})]} - \frac{\mathbb{E}_{\mathbf{y} \sim q}[\nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{y})]}{\mathbb{E}_{\mathbf{y} \sim q}[k(\mathbf{x}, \mathbf{y})]},$$

noting that  $k(\mathbf{x}, \mathbf{y}) \nabla_{\mathbf{x}} \log k(\mathbf{x}, \mathbf{y}) = \nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{y})$ .

## REFERENCES

- Deng, M., Li, H., Li, T., Du, Y., & He, K. (2026, ). *Generative Modeling via Drifting*. <https://arxiv.org/abs/2602.04770>
- Esteban-Casadevall, M., Carrasco-Pollo, J., Welling, M., Meent, J.-W. van de, Bekkers, E. J., & Eijkelboom, F. (2026, ). *Kernel-Gradient Drifting Models*. <https://arxiv.org/abs/2605.10727>
- Gretton, A., Wenliang, L. K., Galashov, A., Thornton, J., Bortoli, V. D., & Doucet, A. (2026, ). *On the Wasserstein Gradient Flow Interpretation of Drifting Models*. <https://arxiv.org/abs/2605.05118>