

Statistical approach “review”

Samy Braik

April 2025

Two paradigms are mentioned in the following sections. The first one is generative models and the second one nonparametric density estimation. These approaches are particularly relevant in our setting because they make few, if any, assumptions on the shape of the data.

On one hand, generative models are the go-to techniques to learn and sample from an unobserved probability distribution. They more or less learn the true distribution, sometimes implicitly, but they are mostly effective at generating new data. On the second hand, density estimation solely focus on the first goal, learning the distribution. Although I argue that we could sample according to a good density estimation.

1 Generative models

All the methods described in this section follow the same framework. They want to link an unknown distribution with density p to a simpler distribution with density q . Either directly like flow methods or up to a certain degree of precision like diffusion models.

1.1 Normalizing flow

Let $X_0 \in \mathbb{R}^d$ distributed according to q a simple distribution, a Gaussian for example, and p a target distribution.

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$, an invertible and differentiable function and set $X_1 := f(X_0)$ such that $X_1 \sim p$.

We can write p as a function of q and f using the change of variable formula

$$p(X_1) = q(f^{-1}(X_1)) \left| \det \frac{\partial f^{-1}}{\partial X_1}(X_1) \right| = q(X_0) \left| \det \frac{\partial f}{\partial X_0}(X_0) \right|^{-1} \quad (1)$$

$$\implies \log p(X_1) = \log q(X_0) - \log \left| \det \frac{\partial f}{\partial X_0}(X_0) \right| \quad (2)$$

Therefore the goal is to learn f_θ , approximation of f , such that $X_1 \simeq f_\theta(X_0)$.

A structure is imposed to f_θ , we define $f_1 \dots f_k$ simple bijective and differentiable transformations, such that

$$f_\theta = f_K \circ f_{K-1} \circ \dots \circ f_2 \circ f_1 \quad (3)$$

where each f_k is chosen such that its Jacobian determinant is easy to compute.

This leads to the following chain of inverse transformations

$$Z_K = X_1, \quad Z_{K-1} = f_K^{-1}(Z_K), \quad Z_{K-2} = f_{K-1}^{-1}(Z_{K-1}), \quad \dots, \quad Z_0 = f_1^{-1}(Z_1), \quad (4)$$

To learn f_θ , we maximize the likelihood of the data, or equivalently minimize the negative log-likelihood.

$$\mathcal{L}_{\text{NF}}(\theta) = -\frac{1}{n} \sum_{i=1}^n \left[\log q(Z_0^{(i)}) + \sum_{k=1}^K \log \left| \det \left(\frac{\partial f_k^{-1}}{\partial z_k}(Z_k^{(i)}) \right) \right| \right], \quad (5)$$

where $Z_0^{(i)} = f_\theta^{-1}(X_i)$, and each $Z_k^{(i)}$ is computed recursively via the inverse transformations.

1.2 Flow

A C^r flow is a time-dependent mapping $\phi : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ characterized by $\phi(t, x) \rightarrow \phi_t(x)$ such that for all $t \in [0, 1]$, ϕ_t is a C^r diffeomorphism in x . We define a flow model by applying a flow ϕ_t to the random variable X_0

$$X_t = \phi_t(X_0), \quad t \in [0, 1], X_0 \sim p \quad (6)$$

Alternatively, we can define a flow using a velocity field $v_t : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ characterized by $v : (t, x) \rightarrow v_t(x)$ via the following ODE

$$\begin{cases} \partial_t \phi_t(x) &= v_t(\phi_t(x)) \\ \phi_0(x) &= x \end{cases} \quad (7)$$

We can derive a probability path as the marginal PDF of a flow model 6 at time t by $X_t \sim p_t$. This PDF is obtained by a push-forward formula

$$p_t(x) = p(\phi^{-1}(x)) |\det \partial_x \phi^{-1}(x)| \quad (8)$$

1.3 Flow matching

Using the notions defined in the previous section, we now introduce the Flow Matching framework. Given a known source distribution q and an unknown target distribution p , we define a probability path p_t interpolating from $p_0 = q$ to $p_1 = p$. We learn a velocity field v_t^θ (a neural network) generating the path p_t and sample according to the learned model by solving the ODE (7).

In order to learn v_t^θ , we minimize the flow matching loss

$$\mathcal{L}_{\text{FM}}(\theta) := \mathbb{E}[\|v_t(X_t) - v_t^\theta(X_t)\|^2] = \mathbb{E}[\|v_t^\theta(X_t) - \dot{X}_t\|^2] + c \quad (9)$$

where $c = \mathbb{E}[\|\dot{X}_t\|^2] - \mathbb{E}[\|v_t(X_t)\|^2]$ constant with respect to θ .

1.4 Diffusion

The general framework of diffusion is divided in two phases. In the forward phase (noising), we start from a random variable distributed according to our target distribution p , gradually add noise until it reaches an easy-to-sample distribution q which is practically always a Gaussian. Then, in the backward phase (denoising) we reverse the process and start from q to get back to p .

We consider $T \in \mathbb{N}^*$ a finite time horizon, a noise schedule $\beta : [0, T] \rightarrow \mathbb{R}_+^*$, assumed to be continuous and non-decreasing, B_t a Brownian motion at time t .

Forward and Backward processes

$$d\vec{X}_t = \frac{-\beta(t)}{2\sigma^2} \vec{X}_t dt + \sqrt{\beta(t)} dB_t, \quad \vec{X}_0 \sim p \quad \text{Forward process} \quad (10)$$

$$d\overleftarrow{X}_t = \left(\frac{\beta(T-t)}{2\sigma^2} \overleftarrow{X}_t + \beta(T-t) \nabla \log p_{T-t}(\overleftarrow{X}_t) \right) dt + \sqrt{\beta(T-t)} d\overleftarrow{B}_t, \quad \overleftarrow{X}_0 \sim p_T \quad \text{Backward process} \quad (11)$$

Since we only noise the random variable until a finite time T , the resulting distribution p_T is not exactly equal to the distribution q , however with a good choice of T and β , we can hope that $p_T \simeq q$. Furthermore,

the backward process allows us to retrieve p but the score ∇p_t is unknown at each time t . To address this problem, denoising score matching is used.

Denoising Score Matching

Let $s : \mathbb{R}^d \rightarrow \mathbb{R}^d$. X a random variable with density p and ε an independent random variable with density g , a centered Gaussian density. Then

$$\mathbb{E}[|\nabla \log p_t(X + \varepsilon) - s(X + \varepsilon)|^2] = c + \mathbb{E}[|\nabla \log g(\varepsilon) - s(X + \varepsilon)|^2] \quad (12)$$

$$= c + \mathbb{E}[|(-\varepsilon/\text{Var}(\varepsilon))g(\varepsilon) - s(X + \varepsilon)|^2] \quad (13)$$

with c a constant not related to s .

With a good architectural choice of the neural network s_θ (data dependent) and noise schedule, we can use the loss to learn the score function and generate new samples from the target distribution using the backward SDE.

2 Nonparametric density estimation

The second approach that could be useful in our situation is nonparametric density estimation. Given an i.i.d. dataset (X_1, \dots, X_n) drawn from a distribution with density f . The goal, like the name suggests, is to estimate f without assuming a specific parametric form. The two most common methods are kernel density estimation and projection estimator.

2.1 Kernel estimator

Consider a kernel function K which is a symmetric density, H a $d \times d$ symmetric and positive definite matrix, $x \in \mathbb{R}^d$, the kernel estimator is defined by

$$\hat{f}_H(x) := \frac{1}{n|H|^{1/2}} \sum_{j=1}^n K\left(H^{-1/2}(x - X_j)\right) = \frac{1}{n} K_H(x - X_j) \quad (14)$$

with $K_H(x) := |H|^{-1/2} K(H^{-1/2}x)$.

A widely used kernel is the Gaussian kernel : $K_H(x) = (2\pi)^{-d/2} |H|^{-1/2} e^{-\frac{1}{2}x^\top H^{-1}x}$

The choice of H is critical since it governs the bias-variance tradeoff and the convergence rate of the estimator. To choose the optimal H few methods could be used like Cross validation or Goldenschluger-Lepski.

2.2 Projection estimator

This method assumes that $f \in L_2(A)$, $A \subset \mathbb{R}^d$.

Let $(\varphi_j)_{j \leq 1}$ be a Hilbert basis of $L_2(A)$ such as Fourier, Legendre or wavelets basis. The projection estimator is defined by

$$\hat{f}_m(x) = \sum_{\|j\| \leq m} \hat{a}_j \varphi_j(x), \quad \hat{a}_j = \frac{1}{n} \sum_{i=1}^n \varphi_j(X_i) \quad (15)$$

Just like the previous case, the choice of m is crucial, and methods like cross validation and penalization help choosing the best model.

References

- [1] Heather Battey and Han Liu. Smooth projected density estimation, 2014.
- [2] Simon Coste. Flow models ii: Score matching techniques, Mar 2025.
- [3] Charlotte Dion-Blanc. Nonparametric density estimation, 2024.
- [4] Tor Fjelde, Emile Mathieu, and Vincent Dutordoir. An introduction to flow matching, January 2024.
- [5] Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky T. Q. Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code, 2024.
- [6] Stanislas Strasman, Antonio Ocello, Claire Boyer, Sylvain Le Corff, and Vincent Lemaire. An analysis of the noise schedule for score-based generative models, 2025.