

## Lecture 13 : Diffusion in Unsupervised Learning

Last week:  $\rightarrow$  Implicit bias of gradient descent: explains why learning is possible in over-parametrised regime, even without reg!  
 $\hookrightarrow$  sometimes this leads to "interesting"/efficient algorithms (non-linear max-margin classif), sometimes not really. (double-descent).

today: Introduction to probabilistic modeling (no labels)  
Key mathematical tool: diffusion.

### Part I: Basic Notions of Probabilistic Modeling.

• Central object: a probability measure  $\mu \in P(\mathcal{R})$ ,  $\mathcal{R}$  is a high-dimensional data domain.

$\rightarrow$  Several situations of practical interest:

(a)  $\mu$  is unknown; we wish to estimate  $\mu$  from iid samples  $x_i \sim \mu$  [Density Estimation; Generative Modeling].

(b)  $\mu$  is known, and we wish to sample from it [Sampling].

(c)  $\mu$  is known up to normalisation;  $\frac{d\mu}{dx}(x) \propto G(x)$ , and we wish to estimate marginals. [Marginalisation]

$\rightarrow$  Several approaches to express probability distributions  $\mu$ .

(a) "Explicit" models that directly parametrise the density.

## ↳ Gibbs / Boltzmann Distributions. (aka Energy-Based Models)

Fix a base measure  $\tau$  in  $\Omega$  (eg Lebesgue) and consider  $f$  of  $f: \Omega \rightarrow \mathbb{R}$ . The Gibbs Measure with energy  $f \in \mathcal{F}$  and inverse temperature  $\beta > 0$  has density

$$\frac{d\mu}{d\tau(x)}(x) = \frac{1}{Z_{\beta f}} e^{-\beta f}, \text{ with}$$

$$Z_{\beta f} = \int_{\Omega} e^{-\beta f(x)} \tau(dx).$$

Examples: Gaussian  $\leftrightarrow f(x) = -(x - m)^T \Sigma^{-1} (x - m)$

Ising Model  $\leftrightarrow f(x) = \sum_{i,j} J_{ij} x_i x_j$   
etc...  $x_i \in \{-1\}$

## ↳ Mixture Models: Consider a measure with density

$$p(x) = \int p(x|h) \cdot dq(h)$$

⊕

$h$ : "latent" variable

$q$ : prior over latent.

$p(x|h)$ : conditional likelihood.

Ex: Gaussian Mixture Models

Variational Autoencoders

:

## ↳ Normalizing Flows. Assume $T_\theta: \Omega \rightarrow \Omega$ is a diffeomorphism. The pushforward of a measure $\nu$ by $T$ is defined as the measure $\mu = T_\theta \nu$ such that

+ A measurable,

$$\mu(A) = \nu(T^{-1}(A))$$

When  $\nu$  has density,  $p(x)$ , then

for any test function  $X$ , we have

$$E_{x \sim \mu} [X(x)] = E_{x \sim \nu} [X(T(x))]$$

$$\int X(x) \cdot q(x) dx \quad \int X(T(x)) \cdot p(x) dx$$

$\Downarrow \quad x = T(x)$

$$\int X(x) \cdot p(T^{-1}(x)) \cdot |DT^{-1}(x)| dx$$

$$\Rightarrow q(x) = p(T^{-1}(x)) \cdot |DT^{-1}(x)|$$

[Normalizing Flows]

[Rezende & Mohamed]

[Tabak, Vanden-Eijnden]

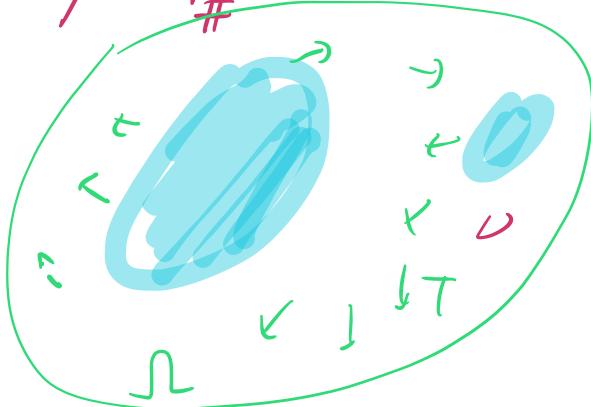
(b) "Implicit" models: define a probability measure indirectly,  $\mu = T_\# \nu$ , where  $T$  is now

a general mapping  $T: \mathbb{D} \rightarrow \mathcal{X}$  (not 1:1, not mapping from  $\mathcal{X}$ ).

Ex: Generative Adv. Nets ( $T$  is given by a NN)

Score-based Diffusion [Song, Ermon, Sohl-Dickstein, ...]  
( $T$  is given by a Reversed - Diffusion probm).

→ Several Criteria for Density Estimation:  $X_i \sim \nu$



drawn iid  $\sim \nu_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  empirical data distribution.

$\mu_\theta$ : one of the probability models from above.

In general, objective is of the form

$$\min_{\theta} \text{"Dist"}(\mu_\theta, \nu_n)$$

Questions:

- Which choice of "Dist"?
- How about "test" error  
"Dist"( $\mu_\theta^n$ ,  $\nu$ )

$\rightarrow$  Kullback - Leibler divergence:  $\text{KL}(\nu || \mu_\theta) = \int \log\left(\frac{d\nu}{d\mu_\theta}\right) d\nu$

$$\text{KL}(\nu || \mu_\theta) = H(\nu; \nu) - H(\mu_\theta; \nu) \quad \text{with}$$

$$H(\mu; \tilde{\mu}) = \int \log(d\mu) \cdot d\mu \quad \text{is the cross-entropy.}$$

Thus minimising KL divergence is equivalent to maximizing cross-entropy.

$\rightarrow H(\mu_\theta; \nu) = E_\nu \left[ \log \frac{d\mu_\theta(x)}{dx} \right]$ , so we can estimate

it from our samples  $x_i$  as

$$\max_{\theta} \frac{1}{n} \sum_{i=1}^n \log \frac{d\mu_\theta}{dx}(x_i) \quad \hookrightarrow \text{Maximum-Likelihood Estimator.}$$

$\rightarrow$  Alternatives to KL?

Integral Probability Metrics (IPMs)

$$D_E(\nu, \mu) = \sup \left\{ |E_\nu f(x) - E_\mu f(x)| \right\}$$

$f \in \mathcal{F}$

$\mathcal{F}$ : family of test functions

Ex 1)  $\mathcal{F}: \{ f: \mathbb{R} \rightarrow [-1, 1] \} \rightsquigarrow \text{TV-distance}$

2)  $\mathcal{F}: \{ f: \Omega \rightarrow \mathbb{R} \text{ 1-Lipschitz} \} \rightsquigarrow 1\text{-Wasserstein distance.}$

In all these cases, we have the same tradeoffs as in SL between approximation, estimation and optimisation errors.

Extra computational challenge: How to efficiently perform training and inference?

Focus on Gibbs/Energy Based Model trained with MLE.

Recall the density of an EBM model is

$$p_\theta(x) = \frac{e^{-\beta f_\theta(x)}}{Z_{\beta, \theta}} = e^{-\beta f_\theta(x) - A(\theta, \beta)}, \text{ with}$$

$$A(\theta, \beta) = \log \int e^{-\beta f_\theta(x)} dx.$$

The MLE becomes

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n -\log p_\theta(x_i) = \underbrace{\frac{\beta}{n} \sum_{i=1}^n f_\theta(x_i)}_{\text{"local" term}} + \underbrace{A(\theta, \beta)}_{\text{"global" term.}} \approx L(\theta)$$

Free energy -  
empirical energy

The gradient  $D_\theta L(\theta)$  becomes

$$D_\theta L(\theta) = \frac{\beta}{n} \sum_{i=1}^n Df_\theta(x_i) - \int \beta Df(\theta) \frac{e^{-\beta f_\theta(x)}}{Z_{\beta, \theta}} dx$$

$p_\theta(x)$

$$= \beta \left( \underbrace{\mathbb{E}_{\nu} [\nabla_{\theta} f_{\theta}]}_{\text{easy to compute.}} - \underbrace{\mathbb{E}_{\mu_{\theta}} [\nabla_{\theta} f]}_{??} \right)$$

→ Natural tool to estimate  $\mathbb{E}_{\mu_{\theta}} [DF]$  is Monte-Carlo:  
 Sample from  $\mu_{\theta}$  and average gradients. But how?  
 ✓  
 MCMC

## Part II: Langevin Diffusion

→ Let's consider the scheme  $X^{(k+1)} = X^{(k)} - \eta \left( \nabla_x f_{\theta}(X^{(k)}) + \sqrt{\frac{1}{\eta}} \cdot \mathcal{N}(0, I) \right)$ .

→ This stochastic discrete process is the discretization of an underlying SDE:

$$dX_t = -\nabla f(X_t) dt + \sqrt{2} \cdot \Gamma dB_t \xrightarrow{\text{Brownian Motion.}} B_t$$

↳ This is a Langevin diffusion.

Convergence of this process?

## Markov Semigroups

→ We can understand the Markov process  $(X_t)_t$  by tracking how measurements  $\mathbb{E}[g(X_t) | X_0 = x]$  evolve over time.

→ Def: For a Markov process  $(X_t)_t$ , its semigroup is a

family  $(P_t)_t$  of operators acting on test functions

$$(P_t g)(x) = \mathbb{E}[g(X_t) \mid X_0 = x]$$

Semigroup properties: (i)  $P_0 = \text{id}$   $P_t \rightarrow \text{id}$  as  $t \rightarrow 0$ .

$$\text{(ii)} \quad P_s P_t = P_{s+t} = P_{t+s} = P_t \cdot P_s$$

→ This semigroup has an associated infinitesimal generator (also a derivative operator):

$$\mathcal{L}g = \lim_{t \rightarrow 0^+} \frac{P_t g - g}{t}$$

→ Generator of Langevin diffusion? ( $\sigma = 1$ )

$$dX_t = -\nabla f(X_t) dt + \sqrt{2} dB_t$$

$$B_t \sim N(0, t \cdot \text{Id})$$

$$X_t = X_0 - \int_0^t \nabla f(X_s) ds + \sqrt{t} B_t$$

$$= X_0 - t \cdot \nabla f(X_0) + \sqrt{t} B_t + o(t)$$

Taylor Expansion of test function  $g$  (Ito Lemma)

$$g(X_t) = g(X_0) + \langle \nabla g(X_0), X_t - X_0 \rangle + \frac{1}{2} (X_t - X_0)^T \nabla^2 g(X_0) (X_t - X_0)$$

$+ o(t)$ .

$$= g(X_0) + \langle \nabla g(X_0), -t \cdot \nabla f(X_0) + \sqrt{t} B_t \rangle + B_t^T \nabla^2 g(X_0) B_t$$

since  $\mathbb{E}[B_t B_t^T]$  is of order  $t$ .  $+ o(t)$

$$\rightarrow \text{Thus } \mathbb{E}[g(X_t) \mid X_0 = x] =$$

$$= g(x) - t \langle \nabla g(x), \nabla f(x) \rangle + t \cdot \text{Tr}(\nabla^2 g(x)) + o(t)$$

$$\therefore \mathcal{L} g(x) = \lim_{t \rightarrow 0^+} \frac{\mathbb{E}[g(X_t) \mid X_0 = x] - g(x)}{t} =$$

$$= - \langle \nabla g(x), \nabla f(x) \rangle + \Delta g(x).$$

$$\rightarrow \text{observe that } \frac{P_{t+h}g - P_tg}{h} = P_t \cdot \frac{P_hg - g}{h} \rightarrow P_t \mathcal{L}$$

"

$$\frac{P_hg - g}{h} \cdot P_t \rightarrow \mathcal{L} \cdot P_t, \text{ so}$$

$$\boxed{\partial_t P_t g = \lim_{h \rightarrow 0^+} \frac{P_{t+h}g - P_tg}{h} = \mathcal{L} P_t g.}$$

↳ this is called the Backwards Kolmogorov Eq.

$$\rightarrow \text{let } (\Pi_t)_t \text{ be the law of the process } X_t. \text{ We can write}$$

$$\mathbb{E} g(X_t) = \int P_t g(x) d\Pi_0(x) = \int g d P_t^* \Pi_0 \quad \begin{matrix} \uparrow \text{formal adjoint of } P_t \\ \end{matrix}$$

$$\text{so } \int g \cdot d\Pi_t \quad \text{so } \Pi_t = P_t^* \Pi_0.$$

. We have

$$\partial_t \int g d(P_t^\bullet \pi_0) = \partial_t \int P_t g d\pi_0 = \int \mathcal{L} P_t g d\pi_0 = \\ = \int P_t \mathcal{L} g d\pi_0 = \int g d(\mathcal{L}^\bullet P_t^\bullet \pi_0) , \text{ so}$$

$$\partial_t P_t^\bullet \pi_0 = \mathcal{L}^\bullet P_t^\bullet \pi_0 = \mathcal{L}^\bullet \pi_t .$$

$$\partial_t \pi_t$$

$$\boxed{\partial_t \pi_t = \mathcal{L}^\bullet \pi_t}$$

Forward-Kolmogorov Eq.  
(Fokker-Plank eq.)

→ let  $g, h : \mathbb{R}^d \rightarrow \mathbb{R}$  smooth.

$$\int \mathcal{L} g \cdot h dx = \int (\Delta g - \langle \nabla g, \nabla f \rangle) \cdot h dx$$

$$= \int g \cdot ("something that depends on h") .$$

$$= \int g (\Delta h + \operatorname{div}(h \cdot \nabla f)) dx .$$

$\uparrow$   
integration  
by parts

$$\text{So } \mathcal{L}^\bullet h = \Delta h + \operatorname{div}(h \cdot \nabla f) .$$

$$\text{So } \partial_t \pi_t = \mathcal{L}^\bullet \pi_t .$$

- "as long / obtaining solution to  $\partial_t \pi = \int \mathcal{L}^\bullet \pi$ "

→ So the equilibrium / stationary solution of  $\partial_t \pi = \alpha L$

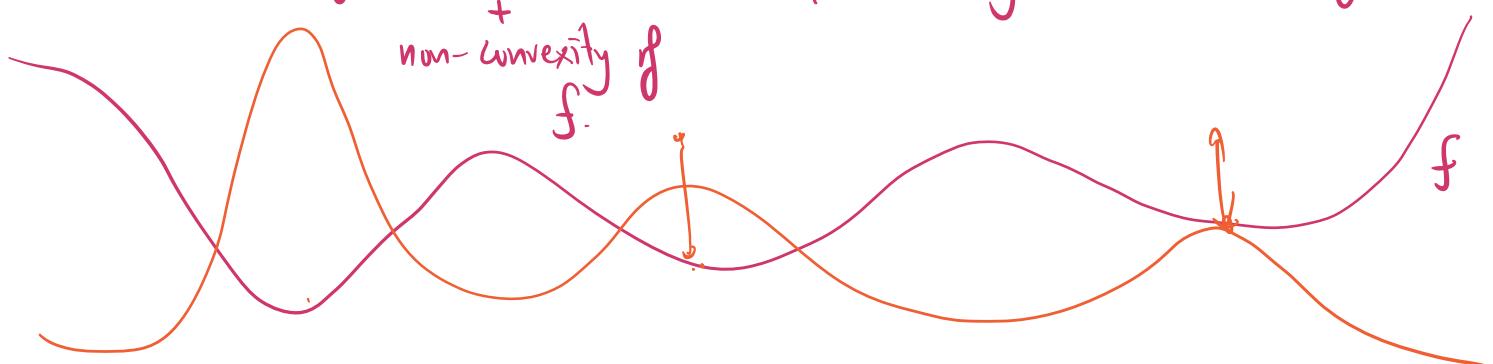
solves  $L^* \pi = 0$

$$0 = \Delta \pi + \operatorname{div}(\pi \nabla f) = \operatorname{div}(\nabla \pi + \pi \nabla f)$$
$$= \operatorname{div}(\pi (\nabla \log \pi + \nabla f))$$

$$\Rightarrow \log \pi = -f + \text{const.} \Rightarrow \pi = \frac{1}{Z} \exp(-f)$$

→ In summary, the unique stationary measure of Langevin diffusion  $dX_t = -\nabla f(X_t) dt + \sqrt{\frac{2}{\beta}} dB_t$  is the Gibbs measure  $\pi_{\beta f}(x) = \frac{1}{Z_B} \exp(-\beta f(x))$ .

→ Although the global convergence to  $\pi_{\beta f}$  holds for very general energies, the rate of convergence is generally cursed by dimension. ↳ presence of "metastability"



→ Rate of convergence of  $\pi_t \rightarrow \pi_*$  is controlled by so-called functional inequalities, measuring how "contractive" the dynamics are in a certain metric

↳ spectral Gap ( $\chi^2$ -divergence)

↳ Log-Sobolev Inequality (KL-divergence)

→ Punch-line / Conclusion: slow rate is unavoidable in high-dim.

→ Sampling hard  $\rightsquigarrow$  EBM learning hard!

Alternatives?

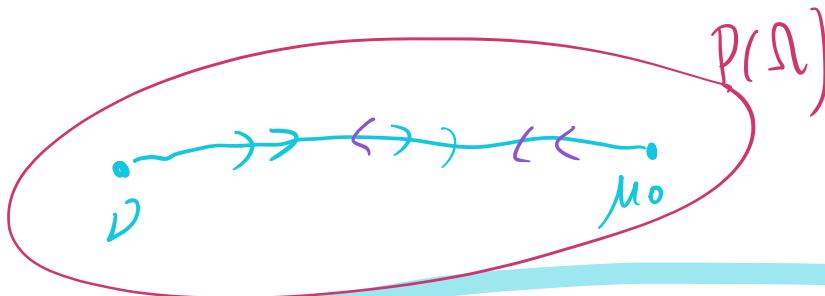
## Score-Based Diffusion

(i) we have seen that in order to sample from  $\mu = \frac{1}{Z} e^{-f}$ ,  
Langevin diffusion only needs access to  $\nabla \log \mu = \nabla f$   
the "Score function".

Idea: Build a diffusion process  $\{X(t)\}_{t=0}^T$  such that:

(i)  $X(0) \sim \nu$  (data distribution) for which we have iid samples.

(ii)  $X(T) \approx \mu_0$  base measure, for which we can sample.



$$dX_t = f(X_t, t) \cdot dt + g(t) d\beta_t$$

Running backwards in time: we sample from  $\mu_0$ , then run  
backwards, until time  $-T$ , gives a sample from  $\nu$ .

→ Reverse Diffusion [Anderson, Follmer]

$$dY_t = \left[ f(x, t) - g(t)^2 \nabla_x \log \pi_t(x) \right] dt + g(t) dB_t.$$

↳ if we can estimate the scores of  $\pi_t$  for  $0 \leq t \leq T$ , then we can go back in time.

→ When  $f$  is linear, the forward diffusion is an Ornstein-Uhlenbeck process, where  $\nabla_x \log \pi_t(x)$  can be efficiently estimated from data!

[Song, Ermon, Sohl-Dickstein et al].

### Questions / Summary:

- Convergence guarantees for score-based wrt Langevin?
- Iterative Sampling schemes  $\hookrightarrow$  Diffusion ↘  
Direct Sampling schemes  $\hookleftarrow$  Transport. [Flows, GAN].