

Generative algorithms/models coming from learning from training data.

e.g. iid samples from p . Try to get some estimate \hat{p} .

Want to sample from \hat{p} rather than evaluate $\hat{p}(x)$
(density estimation)

often have generative model where \hat{p} is implicit and backing out $\hat{p}(x)$ is very difficult.

Suppose \hat{p} is defined implicitly as the probability integral transformation of a density q by a generator function $g: \mathcal{Z} \mapsto \mathcal{X}$.

So for some event A

$$\Pr[A] = \Pr[g^{-1}(A)] = \int_{g^{-1}(A)} q(z) dz = \int_A q(g^{-1}(x)) |\nabla_x g^{-1}(x)| dx$$

$$\text{so } \hat{p}(x) = q(g^{-1}(x)) |\nabla_x g^{-1}(x)|$$

where inverse and Jacobian of g are easy to compute.

Converting a generator to a density estimator is easy.

Is this possible? Yes with normalizing flows, which have nice generators.

In other cases, can still sample from \hat{p} by sampling.

$$\text{If } q, \text{ then evaluate } g(z) \sim \hat{p}$$

Normalizing Flows

want inverse and Jacobian to be simple (and also learnable).

want approach to more complicated $g(\cdot)$ be represented well.

→ use sequence of invertible transformations until desired level of complexity achieved.

Recall if we have invertible mapping $f: \mathbb{R}^d \mapsto \mathbb{R}^d$ with inverse $f^{-1} = g$, and we use this to transform r.v. X with distribution g .

Then $\tilde{X} = f(X)$ has distribution:

$$g(\tilde{x}) = g(x) \left| \det \frac{\partial f^{-1}}{\partial \tilde{x}} \right| = g(x) \left| \det \frac{\partial f}{\partial x} \right|^{-1} \quad (*)$$

where equality is by applying chain rule (inverse function theorem) and is property of Jacobian of invertible functions.

we can construct arbitrarily complex densities by composing several simple maps and successively applying (*), the change of variable formula.

The density $g_k(x)$ obtained by successive transformation of r.v. X_0 with density g_0 through chain of k transformations f_k is

$$X_k = f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1(X_0)$$
$$\ln g_k(X_k) = \ln g_0(X_0) - \sum_{k=1}^k \ln \left| \det \frac{\partial f_k}{\partial X_{k-1}} \right|$$

The path traversed by the r.v. $X_k = f(X_{k-1})$ with initial distribution $g_0(X_0)$ is called the flow, path formed by successive distributions g_k is called normalizing flow.

A property of such transformations is LOTUS (law of the unconscious statistician)

3

$$X \sim p(x) \quad E[h(X)] = \sum_i h(u_i) p_X(u_i)$$

that expectations with respect to the transformed density g_K can be computed without explicitly knowing g_K .

Any expectation $E_{g_K}[h(X)]$ can be written as an expectation

under g_0 :

$$E_{g_K}[h(X)] = E_{g_0}[h(f_K \circ f_{K-1} \circ \dots \circ f_2 \circ f_1(x_0))].$$

which does not require computation of the log det Jacobian terms when $h(X)$ does not depend on g_K .

Infinitesimal flows: natural to consider many functions, limit as we get infinitesimal changes with infinite functions.

Describe how initial density $g_0(x)$ evolves over time

using PDE $\frac{\partial}{\partial t} g_t(x) = \mathcal{J}_t[g_t(x)]$ where \mathcal{J}_t describes

the continuous-time dynamics.

Example: Langevin flow where r.v. X with initial density $g_0(x)$ has densities transformed by Fokker-Planck equation

Discrete Flows

Discrete flows:

determinant of Jacobian corrects for changes to volume in continuous spaces, but volume doesn't change in discrete spaces, so no need to account for it in settings with discrete distributions.

Let X be discrete r.v. and $Y=f(X)$

then pmf of Y is sum over preimage of f :

$$\Pr[Y=y] = \sum_{x \in f^{-1}(y)} \Pr[X=x]$$

for invertible f , $\Pr[Y=y] = \Pr[X=f^{-1}(y)]$

GLOW (how do we do learning)

Let X be high-dimensional random vector with unknown true distribution $X \sim p(x)$.

Collect some iid dataset \mathcal{D} of size N and choose a model $P_\theta(x)$ with parameters θ .

Consider maximum likelihood objective:

$$\mathcal{L}(\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N -\log P_\theta(x^{(i)}) \quad \text{for discrete data.}$$

$$\text{or } \mathcal{J}(\mathcal{D}) \approx \frac{1}{N} \sum_{i=1}^N -\log P_\theta(\tilde{x}^{(i)}) + c$$

where $\tilde{x}^{(i)} = x^{(i)} + u$ with $u \sim \mathcal{U}(0, \alpha)$ and $c = -M \log \alpha$

where α is determined by discretization levels of data and M is data dimension

Optimizing ML objective done through ^{stochastic} gradient descent

when parametric form of p_θ is often a neural network.

In flow-based generation:

$Z \sim p_\theta(z)$ where z is latent variable and $p_\theta(z)$ is often simple like multivariate Gaussian (spherical)

$$X = g_\theta(z).$$

where g_θ is invertible so $z = f_\theta(x) = g_\theta^{-1}(x)$.

since we think of $f = f_1 \circ f_2 \circ \dots \circ f_k$.

$$H_0 = X \xrightarrow{f_1} H_1 \xrightarrow{f_2} \dots \xrightarrow{f_k} Z = H_k.$$

$$\log p_\theta(x) = \log p_\theta(z) + \log \left| \det \frac{\partial z}{\partial x} \right|$$

$$= \log p_\theta(z) + \sum_{k=1}^K \log \left| \det \left(\frac{\partial H_i}{\partial H_{i-1}} \right) \right|$$

scalar quantity "log-det-Jacobian"

can be fairly easy to compute when Jacobian is lower triangular matrix

when Jacobian is triangular:

$$\log \left| \det \left(\frac{\partial H_i}{\partial H_{i-1}} \right) \right| = \text{sum} \left(\log \left| \text{diag} \left(\frac{\partial H_i}{\partial H_{i-1}} \right) \right| \right)$$

sum over all vector elements

elementwise log

diagonal entries.