



SDE  $\rightarrow$  DIFFUSION MODEL

$p(x) \rightarrow$  TARGET

$p_\theta(x)$  Our approximation of  $p(x)$

$$p_\theta(x) = \frac{e^{-f_\theta(x)}}{Z_\theta} \Rightarrow \text{good way to parameterize}$$

because  $-f_\theta(x)$  force everything to be  
more negative.  $Z_\theta$  would have such  
that  $\int p_\theta(x) = 1$

Example:

$$N(\bar{x}/N, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\bar{x}-\mu)^2}{2\sigma^2}}$$

$f_\theta(x) =$  how dense the probability is  
around  $x$ .

Calculating  $Z_\theta$  requires us to integrate  
over the entire space of  $x \Rightarrow$  TOO EXPENSIVE  
so we predict the score

$$\nabla_x \log p_\theta(x) = \nabla_x \log \frac{e^{-f_\theta(x)}}{Z_\theta} =$$

$$= \nabla_x \log e^{-f_\theta(x)} - \cancel{\nabla_x \log z_\theta} = -\nabla_x f_\theta(x)$$

$$= S_\theta(x)$$

SCORE MATCHING:

$$\frac{1}{2} E_{p(x)} \left[ \left\| \nabla_x \log p(x) - S_\theta(x) \right\|_2^2 \right]$$

$\downarrow$   
Score function = where to move  
to increase the probability of a  
data point

But WHO IS  $\nabla_x \log p(x)$ ?

$$\frac{1}{2} E_{p(x)} \left[ \left\| \nabla_x \log(p(x)) - S_\theta(x) \right\|_2^2 \right] =$$

$$= \frac{1}{2} \int p(x) \left( \nabla_x \log(p(x)) - S_\theta(x) \right)^2 dx =$$

$$= \frac{1}{2} \int p(x) \nabla_x^2 \log p(x) + p(x) S_\theta^2(x) - 2 \int p(x) \nabla_x \log p(x) \cdot S_\theta(x) dx$$

$$= \frac{1}{2} \int p(x) \nabla_x^2 \log p(x) + \underbrace{\int p(x) S_\theta^2(x) dx}_{\int p(x) f_\theta(x) S_\theta(x) dx} - \int p(x) \nabla_x \log p(x) \cdot S_\theta(x) dx$$

We know that:  $\nabla_x \log p(x) = \frac{\nabla_x f(x)}{p(x)}$

So we have now:

$$\int p_x \nabla_x \log f(x) s_\theta(x) dx = \int p_x \frac{\nabla_x f(x)}{f(x)} s_\theta(x) dx =$$
$$= \int \nabla_x p(x) s_\theta(x) dx = \text{Integration by parts}$$

$$\int_a^b dv du = uv \Big|_a^b - \int v du \quad \text{so we have:}$$

$$= p(x) s_\theta(x) \Big|_{-\infty}^{+\infty} - \int p_x \nabla_x s_\theta(x) dx =$$

↓  
because  $x \rightarrow +\infty p(x) \rightarrow 0$

=> Going back we have that:

$$\frac{1}{2} E_{p(x)} [ \| \nabla_x \log f(x) - s_\theta(x) \|_2^2 ] =$$

$$= \frac{1}{2} \underbrace{\int p(x) \nabla_x^2 \log(f(x)) dx}_{+} + \frac{1}{2} \int p(x) s_\theta^2(x) dx -$$
$$+ \int p(x) \nabla_x s_\theta(x) dx$$

~~+~~ = there is no  $\theta$  so it is a constant and  
for optimization problem it is 0

$$= \frac{1}{2} \int p(x) s_\theta^2(x) dx + \int p(x) \nabla_x s_\theta(x) dx$$

and  $E_{p(x)} = \int p(x) dx$

$$= \frac{1}{2} E_{p(x)} [S_0^2(x)] + E_{p(x)} [\nabla_x S_0(x)]$$

We see that even if we don't have  $p(x)$

$$\frac{1}{2} E_{p(x)} [\|\nabla_x \log p(x) - S_0(x)\|_2^2] =$$

It doesn't depend on  $p(x)$   $\approx \frac{1}{2} E_{p(x)} [S_0^2(x)] + E_{p(x)} [\nabla_x S_0(x)] + C$

We want to minimize it so if needed, we want  $\nabla_x S_0(x) \rightarrow 0$   $\Rightarrow$  local minimum

while the first term  $E_{p(x)} [S_0^2(x)]$  is 0 at data points

**PROBLEM:**  $\nabla_x S_0(x)$  very expensive

It involves computing the Jacobian especially for images. This has been resolved in:

### Sliced Score Matching: A Scalable Approach to Density and Score Estimation

2<sup>nd</sup> Problem : If the space coverage =>  
 $\Rightarrow$  model is good for data in the  
 data space near the crowded area,  
 but not in others.

To solve this problem we **NOISE** the  
 images (Adding Gaussian noise)

$$\tilde{x} = x + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

$P(x) = P_f(\tilde{x})$  & we solve  
 problem 2.

How much noise depends on  $\sigma$  so:

$$\frac{1}{2} E_{P_f}(\tilde{x}) [\|\nabla_{\tilde{x}} \log P_f(\tilde{x}) - S_0(\tilde{x})\|_2^2]$$

If  $\sigma$  is very large the space is no  
 different

If  $\sigma$  is very small we have 2<sup>nd</sup> problem

$$\frac{1}{2} E_{P_f}(\tilde{x}) [S_0^2(\tilde{x})] + E_{P_f(\tilde{x})} [\nabla_{\tilde{x}} S_0(\tilde{x})]$$

How do we compute  $\nabla_{\tilde{x}} S_0(\tilde{x})$ ?

Pascal Vincent 2010 saw a conference  
Score Matching  $\Leftrightarrow$  Denoising Autoencoders

Train w regular AE but add noise

In the bottleneck  $\Rightarrow$  being able to

separate noise from data it was possible

that the model could learn important  
features. So we start from

$$\frac{1}{2} E_{P_\theta(x)} [\|\nabla_x \log P_\theta(x) - s_\theta(x)\|_2^2]$$

and our goal is:

$$\frac{1}{2} E_{x \sim p(x), \tilde{x} \sim P_\theta(\tilde{x}|x)} [\|s_\theta(\tilde{x}) - \nabla_{\tilde{x}} \log P_\theta(\tilde{x}|x)\|_2^2]$$

So for

$$\frac{1}{2} E_{P_\theta(x)} [\|\nabla_x \log P_\theta(x) - s_\theta(x)\|_2^2] =$$

$$= \frac{1}{2} \int P_\theta(\tilde{x}) (\nabla_{\tilde{x}} \log P_\theta(\tilde{x}) - s_\theta(\tilde{x}))^2 d\tilde{x},$$

$$= \frac{1}{2} \int P_\theta(\tilde{x}) \nabla_{\tilde{x}}^2 \log P_\theta(\tilde{x}) d\tilde{x} +$$

$$+ \frac{1}{2} \int P_\theta(\tilde{x}) S_\theta^2(\tilde{x}) d\tilde{x} +$$

$$-\int p_0(\tilde{x}) \nabla_{\tilde{x}} \log p_0(x) S_0(\tilde{x}) d\tilde{x}$$

$\frac{\nabla_{\tilde{x}} p_0(\tilde{x})}{p_0(\tilde{x})}$

$\downarrow$

$$\int p_0(x) \nabla_{\tilde{x}} p_0(\tilde{x}) S_0(x) d\tilde{x}$$

Use marginalization and so:

$$p_0(\tilde{x}) = \int p(x) p_0(\tilde{x}|x) dx$$

$$= \int \nabla_{\tilde{x}} \left( \int p(x) p_0(\tilde{x}|x) dx \right) S_0(\tilde{x}) d\tilde{x}$$

↓

Cochent integration rule:

$$= \int \left( \int p(x) \nabla_{\tilde{x}} p_0(\tilde{x}|x) dx \right) S_0(\tilde{x}) d\tilde{x}$$

$$= \int \left( \int \int p(x) p_0(\tilde{x}|x) \nabla_{\tilde{x}} \log p_0(\tilde{x}|x) dx \Big| S_0(\tilde{x}) d\tilde{x} \right)$$

= For consistency of integrals we move  
 $s_0(\tilde{x})$

$$= \iint p(x) P_G(\tilde{x}|x) \nabla \log P_G(\tilde{x}|x) s_0(\tilde{x}) d\tilde{x} dx$$

Now we bring everything back:

$$\begin{aligned} & \frac{1}{2} E_{P_G(\tilde{x})} \left[ \| \nabla_{\tilde{x}} \log P_G(\tilde{x}) - s_0(\tilde{x}) \|_2^2 \right] : \\ & = \frac{1}{2} \int_G p(x) \nabla_{\tilde{x}}^2 \log P_G(\tilde{x}) d\tilde{x} + \frac{1}{2} \int_G P_G(\tilde{x}) s_0^2(\tilde{x}) d\tilde{x} \\ & + \iint p(x) P_G(\tilde{x}|x) \nabla \log P_G(\tilde{x}|x) s_0(\tilde{x}) d\tilde{x} dx \\ & = \frac{1}{2} E_{P_G(\tilde{x})} \left[ \| \nabla_{\tilde{x}} \log P_G(\tilde{x}) \|_2^2 \right] + \\ & + \frac{1}{2} E_{P_G(\tilde{x})} \left[ \| s_0(\tilde{x}) \|_2^2 \right] + \\ & - E_{x \sim p(x), \tilde{x} \sim P_G(\tilde{x}|x)} \left[ \nabla_{\tilde{x}} \log (P_G(\tilde{x}|x) \cdot s_0(\tilde{x})) \right] \end{aligned}$$

We open up the maximization

$$= \frac{1}{2} \mathbb{E}_{x \sim p(x), \tilde{x} \sim p_0(\tilde{x}|x)} [\|S_\theta^2(\tilde{x})\|_2^2]$$

and so everything becomes:

$$= \frac{1}{2} \mathbb{E}_{p_0(x)} [\|\nabla_{\tilde{x}} \log p_0(\tilde{x})\|_2^2] +$$

$$+ \frac{1}{2} \mathbb{E}_{x \sim p(x), \tilde{x} \sim p_0(\tilde{x}|x)} [\|S_\theta^2(\tilde{x})\|_2^2 - 2 \nabla_{\tilde{x}} \log p_0(\tilde{x}|x) \cdot S_\theta(\tilde{x})]$$

We still need remove  $\|\nabla_{\tilde{x}} \log p_0(\tilde{x}|x)\|_2^2$ ,

$$\frac{1}{2} \mathbb{E}_{p_0(x)} [\|\nabla_{\tilde{x}} \log p_0(\tilde{x})\|_2^2] +$$

$$+ \frac{1}{2} \mathbb{E}_{x \sim p(x), \tilde{x} \sim p_0(\tilde{x}|x)} [\|S_\theta^2(\tilde{x})\|_2^2 - 2 \nabla_{\tilde{x}} \log p_0(\tilde{x}|x) \cdot S_\theta(\tilde{x})]$$

$$+ \|\nabla_{\tilde{x}} \log p_0(\tilde{x}|x)\|_2^2 - \|\nabla_{\tilde{x}} \log p_0(\tilde{x}|x)\|_2^2$$

$$\frac{1}{2} \mathbb{E}_{x \sim p(x), \tilde{x} \sim p_0(\tilde{x}|x)} [\|S_\theta(\tilde{x}) - \nabla_{\tilde{x}} \log p_0(\tilde{x}|x)\|_2^2]$$

$$= \frac{1}{2} \overline{\mathbb{E}_{P_0(x) \sim P_0}} \left[ \| (\nabla_{\tilde{x}} \log P_0(\tilde{x})) \|_2^2 \right] +$$

$$+ \frac{1}{2} \mathbb{E}_{x \sim P(x), \tilde{x} \sim P_0(\tilde{x}|x)} \left[ \| s_0(x) - \nabla_{\tilde{x}} \log P_0(\tilde{x}|x) \|_2^2 \right] +$$

$$- \| (\nabla_{\tilde{x}} \log P_0(\tilde{x}|x)) \|_2^2$$

$$- \frac{1}{2} \mathbb{E}_{x \sim P(x), \tilde{x} \sim P_0(\tilde{x}|x)} \left[ \| \nabla_{\tilde{x}} \log P_0(\tilde{x}|x) \|_2^2 \right]$$

But we see that ~~-~~ red = do not depend on  $\theta$  & if we want to minimise they are constant red

so our objective becomes:

$$= \mathbb{E}_{x \sim P(x), \tilde{x} \sim P_0(\tilde{x}|x)} \left[ \| s_0(x) - \nabla_{\tilde{x}} \log P_0(\tilde{x}|x) \|_2^2 \right]$$

Remember  $\tilde{x} = x + \epsilon$

$$P_0(\tilde{x}|x) = \frac{1}{(2\pi)^{d/2} \sigma^2} e^{-\frac{1}{2\sigma^2} \| \tilde{x} - x \|^2}$$

So we have that  $\nabla_{\tilde{x}} \log p_f(\tilde{x}|x) =$

$$= \nabla_{\tilde{x}} \log \frac{1}{2\pi \sigma^2} \exp \left( -\frac{1}{2\sigma^2} \|\tilde{x} - x\|^2 \right) =$$

$$= \cancel{\nabla_{\tilde{x}} \log \frac{1}{2\pi \sigma^2}} + \nabla_{\tilde{x}} \log e^{-\frac{1}{2\sigma^2} \|\tilde{x} - x\|^2}$$

$$= -\frac{1}{2\sigma^2} \cancel{2\|\tilde{x} - x\|} = -\frac{\tilde{x} - x}{\sigma^2} =$$

$$= \frac{\tilde{x} - x}{\sigma^2}$$

So we have that:

$$\frac{1}{2} \mathbb{E}_{x \sim p(x), \tilde{x} \sim p_f(\tilde{x}|x)} \left[ \|s_\theta(\tilde{x}) - \frac{1}{\sigma^2} (x - \tilde{x})\|_2^2 \right]$$

Remember that  $\tilde{x} = x + \epsilon$

$$= \frac{1}{2} \mathbb{E}_{x \sim p(x), \tilde{x} \sim p_f(\tilde{x}|x)} \left[ \|s_\theta(x) + \frac{\epsilon}{\sigma^2}\|_2^2 \right]$$

So we can see the our model so predicts the negative of the noise that we added. The negative of the noise points in the direction of the data manifold where the data likely least grows

## HOW DO WE GENERATE NEW DATA?

If we do just see there is an issue because we overlooked the above said fail so we do it step by step with small step

$$\tilde{x}_{i+1} \leftarrow \tilde{x}_i + \alpha \cdot s_0(\tilde{x}_i)$$

The size of the step is determined by  $\alpha$

The more points to the center because the likelihood is high because

Now's where we have higher density  
of dots

There is an easy fix, it is called  
**LANGEVIAN DYNAMICS**

$$\tilde{x}_{i+L} \leftarrow \tilde{x}_i + \alpha s_0(\tilde{x}_i) + \sqrt{\alpha} \cdot \epsilon$$

we add this  $\epsilon \sim$  Random  
Gaussian noise

$\sqrt{\alpha}$  very small

In this way we have nice artifacts

that do not collapse

---

Considering 6 damping levels,

so the model can be distributed  
over all over the space

so we have  $\tilde{x} = x + \epsilon \sim N(0, \sigma^2 I)$

so  $\sigma$  is not fixed anymore

If you increase  $t \rightarrow +\infty$  the noise perturbation becomes stochastic

so we can use SDEs to model stochastic processes.

$$dx = f(x, t) dt + g(t) dw$$

drift      diffusion      Wiener process  
infinitesimal changes in time      infinitesimal changes in noise

drift is fixed and describes if and how we move  $x$  in a deterministic way

diffusion describes the influence of the stochasticity over time

the more motility we have

$$\tilde{x} = x + G \quad \epsilon \sim \mathcal{N}(0, \sigma_t^2 I)$$

So change of  $x$  is only influenced by stochastic term so  $f(x, t) = 0$

$$dx = g(t) dw \quad \text{to per } t \rightarrow +\infty$$

$$\sigma_t \geq \sigma(t)$$

$$\downarrow \\ G \in [0, 1]$$

N.B.

IN DDPM  $f(x, t) \neq 0$

Thus SDE corresponds to the process  
of missing our data FORWARD SDE

$$dx = g(t) dw$$

REVERSE SDE (WITH LANGEVIAN DYNAMICS)

$\downarrow$   
Sampling

$$\text{If FORWARD } dx = f(x, t) dt + g(t) dw$$

The REVERSE SDE is:

$$dx = [f(x_t, t) - \beta^2(t) \nabla_x \log p_{\theta}(x)] dt + \sigma(t) dw$$

Score function

So to reverse we can use the  
SCORING FUNCTION

## CONNECTION TO DDPM

Score matching  $\tilde{x} = x + \epsilon \in \mathcal{N}(0, \beta_t^2)$

DDPM  $x_t = \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon \in \mathcal{N}(0, I)$

So we see DRAFT is now zero

If we compute the reverse SDE for  
DDPM we have near

$$dx = [f(x_t, t) - \beta^2(t) \nabla_x \log p_{\theta}(x)] dt + \sigma(t) dw$$

$$dx = x_t - x_{t-1}$$

$$= \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon - x_{t-1}$$

$\downarrow$  Taylor expansion =  $1 - \frac{1}{2} \beta_t$

$$= \left(1 - \frac{1}{2} \beta_t\right) x_{t-1} + \sqrt{\beta_t} \epsilon - x_{t-1} :$$

~~$$= x_{t-1} - \frac{1}{2} \beta_t x_{t-1} + \sqrt{\beta_t} \epsilon - x_{t-1} :$$~~

~~$$= - \frac{1}{2} \beta_t x_{t-1} + \sqrt{\beta_t} \epsilon$$~~

So e)  $x_t - x_{t-1} \rightarrow 0$ ,  $t \rightarrow \infty$

$$\boxed{dx = -\frac{1}{2} \beta_t x dt + \sqrt{\beta_t} dw}$$

FORWARD SIDE

REVERSE SDE  $f(x, \epsilon) = -\frac{1}{2} \beta_t x$

$$f(t) = \sqrt{\beta_t}$$

We apply Ito formula

$$dx = \left[ -\frac{1}{2} \beta(t) x - \beta(t) \sqrt{\log P_0(x)} dt + \right. \\ \left. + \sqrt{\beta(t)} dw \right]$$

If we discretize  $t$ :



### DDPM SAMPLER

$$x_{t-1} = \frac{1}{\sqrt{1-\beta_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\tilde{\alpha}_t}} s_0(x_t, t) \right) \\ + \sqrt{\beta_t} z$$

We discretize  $t$  with the

Euler Maruyama Method