

There is a quantity that captures the Signal-to-Noise ratio indeed, and it was named the same. It is related to the noise coefficients that we are free to choose in the $q(X_t|X_0)$ samples. To recall, $X_t = \alpha_t X_0 + \Sigma_t \epsilon$.

Previously we had $\alpha_t = \sqrt{\lambda_t}$ and $\Sigma_t = \sqrt{1-\lambda_t}$, where the λ_t is a sequence that approaches 0 as $t \rightarrow \infty$ or approximately 0 when $t \rightarrow T$.

What we require about α_t and Σ_t is that $\alpha_t \rightarrow 0$, $\Sigma_t \rightarrow 1$ as $t \rightarrow T$ because we have to approach a standard normal $N(0, I)$.

So, this Signal-to-Noise ratio, denoted

$$\text{SNR}(t) = \frac{\alpha_t^2}{\Sigma_t^2}. \quad \text{It is a measure of noise.}$$

Although, if $\text{SNR}(T) = \frac{\alpha_T^2}{\Sigma_T^2} = 0$ because $\alpha_t \rightarrow 0$, ~~we~~ it actually tells us, in a sense,

that the initial image will

more much from the initial image is still present at step t .

Also, another quantity of interest "the intensity" is $\lambda = \log \text{SNR}$. Given that t is a random variable (uniform as in DDPM), λ above is a change of coordinates in the PDF'n.

$$p(\lambda) = p(t) \left| \frac{dt}{d\lambda} \right|$$

Assuming λ is decreasing with t (which is natural to assume under the limits above):

$$p(\lambda) = -\frac{dt}{d\lambda}.$$

As an example, the cosine schedule is $a_t = \cos(\frac{\pi t}{2})$, $\nabla_t = \sin(\frac{\pi t}{2})$ ($\forall t \in [0, 1]$, otherwise $\frac{\pi t}{2T}$, for $T = 1000$ or other). So, $\lambda = -2 \ln \tan \frac{\pi}{2}$, $t = \frac{2}{\pi} \arctan e^{-\lambda/2}$. Thus, $p(\lambda) = -\frac{dt}{d\lambda}$

$$= \frac{1}{2\pi} \text{sech} \frac{\lambda}{2}.$$

$$\text{Also, } t = 1 - \int_{-\infty}^{\infty} p(\lambda) d\lambda = P(\lambda)$$

$$\Rightarrow \lambda = P^{-1}(t).$$

We can also obtain the noise schedule λ by the inverse P^{-1} .

Under a general SNR we still have to look at how the ELBO changes.

$$\text{Let } \mathcal{L}(x_0) = \sum_{i=1}^T \mathbb{E}_{q(x_i | x_0)} D_{KL}[q(x_{i-1} | x_i, x_0) p(x_{i-1} | x_i, x_0)].$$

There is a similarity between this and an integral if the steps are small and $T \rightarrow \infty$, transforming the sum into an integral.

For it, we will adopt a different notation by discretizing the $[0, 1]$ interval into steps

$$n(i) = \frac{i-1}{T} := \alpha \quad \text{and} \quad t(i) = \frac{i}{T} := t$$

$$p(x_n | x_t) = q(x_n | x_t, \hat{x}_\theta(x_t; t))$$

$$= \mathcal{N}(x_n; \mu_q(x_n, x_t; \alpha, t), \Sigma_q(\alpha, t) \mathbf{I})$$

$$\text{and } p(x_n | x_t) = \mathcal{N}(x_n; \mu_\theta(x_t; \alpha, t), \Sigma_q(\alpha, t) \mathbf{I}).$$

Recall this is the setup that we had in DDPM.

$$\mu_q = \frac{\frac{\alpha_t}{\alpha_n} \Sigma_n^2}{\Sigma_t^2} x_t + \frac{\alpha_n}{\Sigma_t^2} \left(\Sigma_t^2 - \frac{\alpha_t^2}{\alpha_n^2} \Sigma_n^2 \right) x_0$$

$$\mu_\theta = \frac{\frac{\alpha_t}{\alpha_n} \Sigma_n^2}{\Sigma_t^2} x_t + \frac{\alpha_n}{\Sigma_t^2} \left(\Sigma_t^2 - \frac{\alpha_t^2}{\alpha_n^2} \Sigma_n^2 \right) \hat{x}_\theta$$

$$\Sigma_q^2 = \left(\Sigma_t^2 - \frac{\alpha_t^2}{\alpha_n^2} \Sigma_n^2 \right) \frac{\Sigma_n^2}{\alpha_t^2}$$

After computing the KL-div, we end up with:

$$D_{KL} = \frac{1}{2} (\text{SNR}(0) - \text{SNR}(t)) \|x - \hat{x}_{\theta}(x; t)\|_2^2$$

With $x_t = \alpha_t x_0 + \sqrt{1 - \alpha_t} \varepsilon$, $\varepsilon \sim \mathcal{N}(0, I)$

$$\mathcal{L}(x_0) = \frac{1}{2} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, I), i \sim \{1, \dots, T\}} [(\text{SNR}(0) - \text{SNR}(t)) \|x - \hat{x}_{\theta}(x_t; t)\|_2^2]$$

When $T \rightarrow \infty$ we see here the sum converges to an integral

$$\mathcal{L}_{\infty}(x_0) = -\frac{1}{2} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, I)} \int_0^1 \text{SNR}'(t) \|x - \hat{x}_t\|_2^2 dt$$

The new variational lower bound indicates an equivalence under scaling which might be used as a new loss, where during training a range of noisy samples are prioritized.

$$\text{SNR}'(t) = (e^{\lambda})' = \frac{d\lambda}{dt} \text{SNR} \quad \text{or} \quad \frac{1}{\frac{dt}{d\lambda}} \text{SNR}$$

$= \frac{w(\lambda)}{\frac{dt}{d\lambda}}$, if we compactify the above expression

$p(\lambda)$ s.t. everything except $p(\lambda)$ is
some general weighting $w(\lambda)$!

Our new weighted loss becomes

$$\mathcal{L}_w(x_0) = -\frac{1}{2} \mathbb{E}_{t \sim U(0,1), \epsilon \sim N(0, I)}$$

$$\left[\frac{w(\lambda_t)}{p(\lambda_t)} \quad \|\hat{\epsilon}_\theta - \epsilon\|_2^2 \right]$$

The minus comes from the fact we want to
minimize, and noise time is sampled uniformly
in $(0,1)$, although for practical reasons can
be converted to $[0,1]$ in implementation.
Additionally, if we replace $\frac{d\lambda}{dt}$ in $\mathcal{L}_\infty(x_0)$
we see now the integral no longer depends
on $p(\lambda)$ and only $w(\lambda)$. So, the integral
is independent of the noise schedule (remember,
 $w(\lambda)$ we want to make it arbitrary,
although from our deduction looks like
Now the bounds change

$w(\lambda) \sim \text{SNR}$. Hence, the

$\int_{\lambda_{\min}}^{\lambda_{\max}}$ as long as two noise schedulers $p(\lambda)$ define same $\lambda_{\max}, \lambda_{\min}$, the integral is the same (and $w(\lambda)$ is unchanged).

The last lines can be ignored, since we aren't computing the time integral, but only sampling as a Monte Carlo approximation. Due to sampling, the gradients and variance of the expectation depends on $p(\lambda)$ or noise scheduler. Hence, there is a reason to search for different noise schedulers to probably optimise sampling.

The final loss function becomes:

$$\mathcal{L}_w(\theta) = \frac{1}{2} \mathbb{E}_{x \sim \mathcal{D}, \epsilon \sim \mathcal{N}(0, I), \lambda \sim p(\lambda)} \left[\frac{w(\lambda)}{p(\lambda)} \| \epsilon_{\theta}(x_{\lambda}; \lambda) - \epsilon \|_2^2 \right].$$

The noise schedule that was tested to achieve the best sampling quality is Laplace:

$$p(\lambda) = \frac{e^{-\frac{|\lambda - \mu|}{b}}}{2b}$$

$$\lambda(t) = \mu - b \operatorname{sgn}(0,5 - t) \ln(1 - 2|t - 0,5|)$$