

$\alpha$ -STABLE DIFFUSIONS:

Extreme Diffusion.

The main goal of this notebook lies in developing  $k$ -stable diffusions. It is based on the following observation: a lot of what makes diffusion models tractable is that the normal distribution is stable.

That is, (from Wikipedia):

Let  $X_1$  and  $X_2$  be independent realizations of a random variable  $X$ . Then  $X$  is said to be stable if for any constants  $a > 0$  and  $b > 0$  the random variable  $aX_1 + bX_2$  has the same distribution as  $cX + d$  for some constants  $c > 0$  and  $d$ . The distribution is said to be strictly stable if it holds for  $d = 0$ .

The normal, Cauchy and Lévy distributions all have the same property.

The probability density function for a general stable distribution cannot be written analytically, (Why?). A random variable that is stable has the following characteristic function, in general,

$$\varphi(t; \alpha, \beta, \gamma, \mu) = \exp(i\mu t - |\gamma t|^\alpha (1 - i\beta \operatorname{sign}(t) \phi))$$

$$\phi = \begin{cases} \tan\left(\frac{\pi\alpha}{2}\right) & \alpha \neq 1 \\ -\frac{2}{\pi} \log|t| & \alpha = 1 \end{cases}$$

$\alpha \in (0, 2]$  shape - asymmetry  
 $\beta \in (-1, 1)$  concentration  
 $\gamma \in (0, \infty)$  scale  
 $\mu \in (-\infty, \infty)$  location

## Properties

- All stable distributions are infinitely divisible. iid r.v.'s.
- With the exception of the normal distribution ( $\alpha=2$ ), stable distributions are leptokurtic and heavy-tailed distributions
- Closure under convolution.
  - ↳ Can be seen from multiplying chrs.
  - ↳ function.

→ if can be expressed  
as the sum of an  
arbitrary # of

$\alpha$ -stable distributions generalize the C.T!

Generalized C.T.

A non-degenerate random variable  $Z$  is  $\alpha$ -stable for some  $0 < \alpha \leq 2$  if and only if there is an independent, identically distributed sequence of random variables  $X_1, X_2, \dots$  and constants  $a_n > 0$ ,  $b_n \in \mathbb{R}$  with

$$a_n(X_1 + \dots + X_n) - b_n \xrightarrow{d} Z$$

## Cauchy Distribution.

A relatively small number of  $\kappa$ -stable distributions have a representation in terms of simple functions. Apart from the normal distribution, ( $\kappa=2$ ), the Cauchy distribution is the most studied. Here are a few properties (Wikipedia):

1. If  $X \sim \text{Cauchy}(x_0, \gamma)$  then  $kX + l \sim \text{Cauchy}(x_0 k + l, \gamma/k)$ .
2. If  $X \sim \text{Cauchy}(x_0, \gamma_0)$  and  $Y \sim \text{Cauchy}(x_1, \gamma_1)$  are independent, then  $X + Y \sim \text{Cauchy}(x_0 + x_1, \gamma_0 + \gamma_1)$  and  $X - Y \sim \text{Cauchy}(x_0 - x_1, \gamma_0 + \gamma_1)$ .
3. If  $X \sim \text{Cauchy}(0, \gamma)$  then  $\frac{1}{X} \sim \text{Cauchy}(0, \frac{1}{\gamma})$ .
4. Has univariate pdf

$$f(x) = \frac{1}{\pi \gamma \left[ 1 + \left( \frac{x-x_0}{\gamma} \right)^2 \right]}$$

$x \sim \text{Cauchy}(x_0, \gamma)$

Has multivariate pdf:

$$f(\mathbf{x}; \mu, \Sigma, \kappa) = \frac{\Gamma\left(\frac{1+\kappa}{2}\right)}{\Gamma\left(\frac{1}{2}\right) \pi^{\frac{k}{2}} |\Sigma|^{\frac{1}{2}} \left[ 1 + (\mathbf{x}-\mu)^T \Sigma^{-1} (\mathbf{x}-\mu) \right]^{\frac{1+\kappa}{2}}}$$

Note:  $\Sigma = I$  does not imply that the r.v.'s are independent.

5. Has univariate chf:

$$\psi_X(t) = E[e^{itX}] = e^{ix_0 t - \gamma |t|}$$

and multivariate chf

$$\psi_{\mathbf{X}}(\mathbf{t}) = e^{iX_0(t) - \gamma(t)}$$

real functions of degree one s.t.  $X_0(at) = aX_0(t)$   
 $\gamma(at) = |a|\gamma(t)$ .

6. The KL-divergence between two Cauchy distributions has a closed form

$$KL(P_{X_0, \gamma_1}; P_{X_0, \gamma_2}) = \log \frac{(x_1 + y_2)^2 + (x_{0,1} - x_{0,2})^2}{4\gamma_1 \gamma_2}$$

## Motivation of Question:

DDPM's algorithm for training is:

1. Repeat
2.  $x_0 \sim q(x_0)$   $\mapsto$  Sample from data.
3.  $t \sim \text{Unif}\{\dots, T\}$ .  $\mapsto$   $T$  is # of timesteps.
4.  $\epsilon \sim N(0, I)$ .

5 Take gradient step on.

$$\nabla_{\theta} \|\epsilon - \epsilon_0(x_t, t)\|^2 \mapsto \nabla_{\theta} \|\epsilon - \epsilon_0(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, t)\|^2$$

6. Until converged

A crucial step, (5), ~~exist~~ is dependent upon the stability of the normal, because for some schedule of  $\beta_t$ 's,

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon$$

Does the Cauchy distribution, with its stability and existing KL divergence determine a viable diffusion process with different or desirable (extremes?) properties.

from the review ~~and~~ there are a few things that worry me:

- (1) - stability might not be enough? Maybe conjugacy is related to stability, right?
- (2) - too much needed? Maybe too much?
- (3) - will it be enough to exchange  $\mu$  and  $\sigma$  to  $x_0$ , and  $\lambda$ ?  
 Ho's DDPM claims that for  $\beta_t$  small, the reverse process is also Gaussian. He cites Sollich-Pickerstein (2015) for this claim, who in turn cites Feller. (an old article). How does this claim hold for the Cauchy case?

With all that in mind, the simplest litmus test for this idea might be assuming univariate distributions (bivariate top), and run with it.

- A. Gaussian Diffusion Process gradually transforms data  $x$  into random noise by adding increasing amounts of Gaussian noise at each timestep  $t$ .

O

Bishop's New book

has a section on diffusion, and although I find it lacking, it does have one stellar problem, that I hope will help: problem 20.7, where he asks you to show that the reverse  $q(z_t | z_{t+1})$  is approximately Gaussian when  $\beta_t$  is small. If I can reproduce this, then I might be able to reproduce it for the Cauchy case.

Anyway, Bishop's diffusion:

Forward Encoder:

$$z_1 = \sqrt{1-\beta_1}x + \sqrt{\beta_1}\epsilon_1, \quad \epsilon_1 \sim N(\epsilon_1 | 0, I).$$

then  $q(z_1 | x) = N(z_1 | \sqrt{1-\beta_1}x, \beta_1 I).$

More generally,  $z_t = \sqrt{1-\beta_t}z_{t-1} + \sqrt{\beta_t}\epsilon_t, \quad \epsilon_t \sim N(\epsilon_t | 0, I).$

$$q(z_t | z_{t-1}) = N(z_t | \sqrt{1-\beta_t}z_{t-1}, \beta_t I).$$

$$\beta_1 < \beta_2 < \dots < \beta_T.$$

## Diffusion Kernel.

$$q(z_1, \dots, z_t | X) = q(z_1 | X) \prod_{\tau=2}^t q(z_\tau | z_{\tau-1}).$$

One can marginalize over intermediate variables and find that:

$$q(z_t | X) = N(z_t | \sqrt{\kappa_t} X, (1-\alpha_t) I)$$

$$\alpha_t = \prod_{\tau=1}^t (1 - \beta_\tau),$$

or  $z_t = \sqrt{\alpha_t} X + \sqrt{1-\alpha_t} \epsilon_t, \quad \epsilon_t \sim N(\epsilon_t | 0, I).$

and as  $T \rightarrow \infty$ , we have

$$q(z_T | X) = N(z_T | 0, I).$$

and  $\underbrace{q(z_T)}_{\text{Indep. of } X} = N(z_T | 0, I)$

## Conditional distribution.

We are trying to unlearn to undo the noise process, and so, we want  $q(z_t | z_{t-1})$ :

$$q(z_t | z_{t-1}) = \frac{q(z_t | z_{t-1}) q(z_{t-1})}{q(z_t)}$$

and

$$q(z_{t-1}) = \int q(z_{t-1} | X) p(x) dx.$$

where  $q(z_{t-1} | X) = N(z_{t-1} | \sqrt{\kappa_{t-1}} X, (1-\alpha_{t-1}) I)$

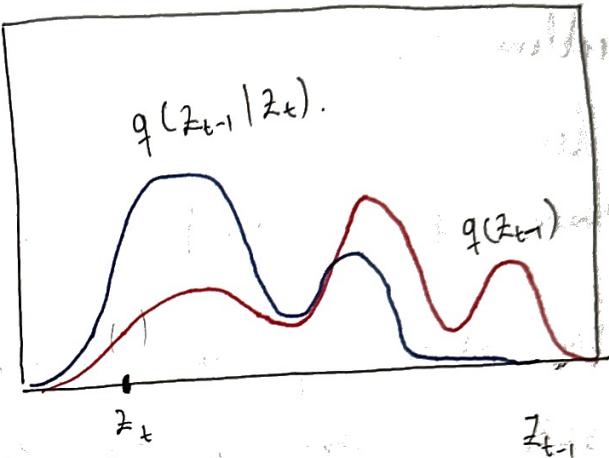
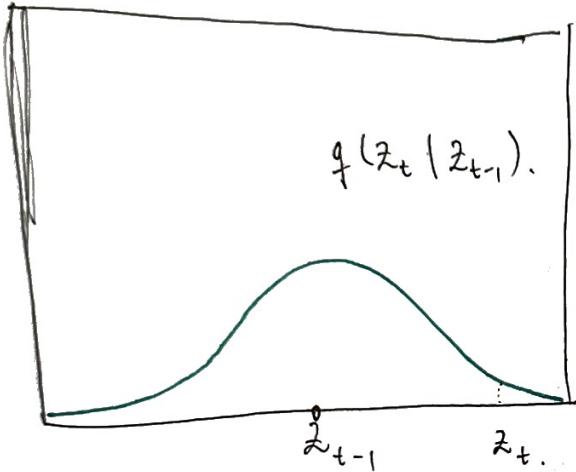
Hence  $q(z_{t-1} | z_t)$  is intractable, so we condition on  $X$ .

The proof is left to myself in the future  $q(z_{t-1} | z_t, X) = \frac{q(z_t | z_{t-1}, X) q(z_{t-1} | X)}{q(z_t | X)}$

$$\hookrightarrow q(z_{t-1} | z_t, X) = N(z_{t-1} | \frac{(1-\alpha_{t-1})\sqrt{1-\beta_t} z_t + \sqrt{\kappa_{t-1}} \beta_t X}{1-\alpha_t}, \frac{\beta_t (1-\alpha_{t-1})}{1-\alpha_t} I)$$

## Reverse Decoder.

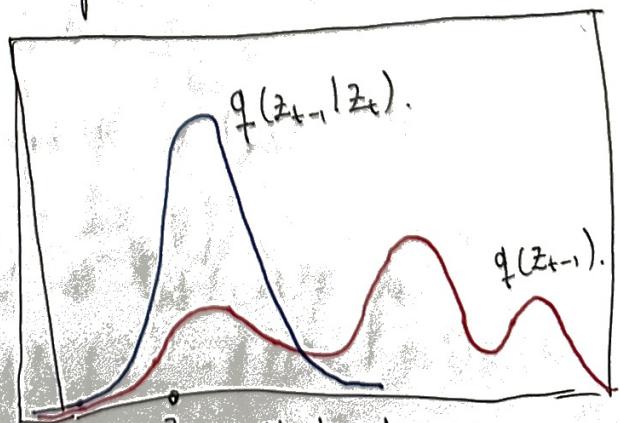
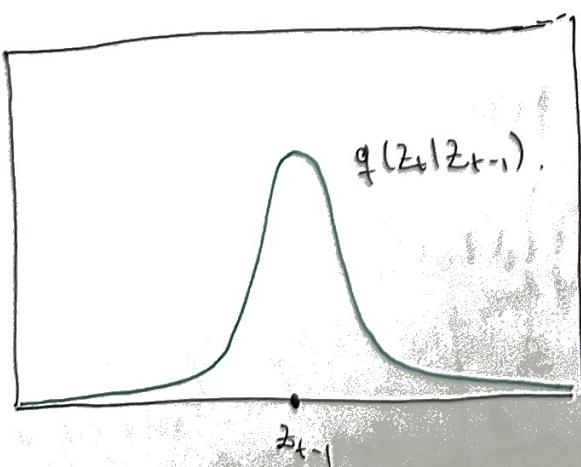
But we really want  $q(z_{t-1}|z_t)$  because if we are doing generative modeling beginning from  $z_t$ , we can not assume the existence of  $x$ .  
 (knowledge)



$$q(z_t | z_{t-1}) = N(z_t | \sqrt{1-\beta_t} z_{t-1}, \beta_t I).$$

$$q(z_{t-1} | z_t) = \frac{q(z_t | z_{t-1}) q(z_{t-1})}{q(z_t)}.$$

If  $q(z_t | z_{t-1})$  is a sufficiently narrow Gaussian then  $q(z_{t-1})$  will vary only a small amount over the region in which  $q(z_t | z_{t-1})$  has significant mass, and hence  $q(z_{t-1} | z_t)$  will be approximately Gaussian.



Maybe the same idea works for Cauchy diffusion!

More formally,  $q(z_{t-1}|z_t)$  will be approximately Gaussian by making a Taylor expansion of  $\ln q(z_{t-1}|z_t)$  around the point  $z_t$  as a function of  $z_{t-1}$ . This also shows that for small variance, the reverse distribution  $q(z_t|z_{t-1})$  will have a covariance that is close to  $\beta_t I$ .

Problem 20.7.

Consider the inverse conditional  $q(z_{t-1}|z_t) = \frac{q(z_t|z_{t-1})q(z_{t-1})}{q(z_t)}$   
where  $q(z_t|z_{t-1}) = N(z_t; \sqrt{1-\beta_t} z_{t-1} + \beta_t I)$

By taking log and making a Taylor expansion of  $q(z_{t-1})$  centred on  $z_t$ , show that, for small  $\beta_t$ ,  $q(z_{t-1}|z_t)$  is approximately Gaussian with mean  $z_t$  and covariance  $\beta_t I$ .

$$\log q(z_{t-1}|z_t) = \log q(z_t|z_{t-1}) + \log q(z_{t-1}) - \log q(z_t)$$

$$\log q(z_{t-1}|z_t) \approx \log q(z_t|z_{t-1}) + \log q(z_{t-1}) - (\log q(z_t))$$

$$q(z_{t-1}|z_t) \propto q(z_t|z_{t-1}) q(z_{t-1})$$

$$\log q(z_{t-1}|z_t) \propto \log q(z_t|z_{t-1}) + \log q(z_{t-1}).$$

$$\log \frac{q(z_{t+1} | z_t)}{q(z_t | z_t)} \approx \log q_{z_{t+1}}(z_t) + (z_{t+1} - z_t)^T D q_{z_{t+1}}(z_t)$$

$$\log q(z_{t+1} | z_t) \approx -\frac{1}{2\beta_t} (z_t - \sqrt{1-\beta_t} z_{t-1})^T (z_t - \sqrt{1-\beta_t} z_{t-1})$$

$$+ \log \frac{q_{z_{t+1}}(z_t)}{\log \frac{2\pi}{\text{constants w.r.t. } z_{t+1}}} + (z_{t+1} - z_t)^T D q_{z_{t+1}}(z_t)$$

$$= -\frac{1}{2\beta_t} (z_t - \sqrt{1-\beta_t} z_{t-1})^T (z_t - \sqrt{1-\beta_t} z_{t-1})$$

$$+ (z_{t+1} - z_t)^T D q_{z_{t+1}}(z_t) + \text{const.}$$

$$= -\frac{1}{2\beta_t} ((1-\beta_t) z_{t-1}^T z_{t-1} - 2\sqrt{1-\beta_t} z_t^T z_{t-1}) + D_t^T z_{t-1} + \text{const}$$

$$= -\frac{1}{2\beta_t} ((1-\beta_t) z_{t-1}^T z_{t-1} - 2(\sqrt{1-\beta_t} z_t + \beta_t D)^T z_{t-1}) + \text{const}$$

$$= -\frac{1}{2\beta_t} \frac{1}{1-\beta_t} (z_{t-1}^T z_{t-1} - 2(\frac{1}{\sqrt{1-\beta_t}} z_t + \frac{\beta_t}{1-\beta_t} D)^T z_{t-1}) + \text{const}$$

$$= -\frac{1}{2\beta_t} \frac{1}{1-\beta_t} (z_{t-1} - (\frac{1}{\sqrt{1-\beta_t}} z_t + \frac{\beta_t}{1-\beta_t} D))^T (z_{t-1} - (\frac{1}{\sqrt{1-\beta_t}} z_t + \frac{\beta_t}{1-\beta_t} D))$$

Expansions at  $\beta=0$ .

$$\sqrt{1-x} \approx 1 + \frac{x}{2} + \frac{3x^2}{8}$$

$$\frac{x}{1-x} \approx x + x^2 + x^3$$

$$\frac{1}{1-x} \approx 1 + x + x^2 + \dots$$

small  $\beta_t$

$$\approx -\frac{1}{2\beta_t} (1) (z_{t-1} - ((1)z_t + \beta_t D))^T (z_{t-1} - (z_t + \beta_t D))$$

$$\xrightarrow{\beta_t \rightarrow 0} -\frac{1}{2\beta_t} (z_{t-1} - z_t)^T (z_{t-1} - z_t) \sim N(z_t, \beta_t)$$

Can we apply this to a Cauchy diffusion?

Let's define it:

Forward Euler:

Consider  $x$  given, or in the univariate setting.

$$z_1 = \sqrt{1-\beta_1} x + \sqrt{\beta_1} \epsilon_1, \quad \epsilon_1 \sim \text{Cauchy}(\epsilon_1 | 0, 1). \quad (1)$$

$$\Rightarrow q(z_1 | x) = \text{Cauchy}(z_1 | \sqrt{1-\beta_1} x, \sqrt{\beta_1} I)?$$

From property (1) of the previous page on the Cauchy distribution:

1. If  $\epsilon \sim \text{Cauchy}(\epsilon_0, \gamma)$  then  $k\epsilon + l \sim \text{Cauchy}(\epsilon_0 k + l, \gamma^2 k)$ .

Hence, if  $\epsilon_1 \sim \text{Cauchy}(\epsilon_1 | 0, 1)$ , then:

$$z_1 | x \sim \text{Cauchy}(\sqrt{1-\beta_1} x, \sqrt{\beta_1}). \quad (2)$$

Each successive <sup>point</sup> ~~stage~~ is given by:

$$(3) \quad z_t = \sqrt{1-\beta_t} z_{t-1} + \sqrt{\beta_t} \epsilon_t \quad \epsilon_t \sim \text{Cauchy}(\epsilon_t | 0, 1)$$

Where again, by the Markov property and 1.

$$q(z_t | z_{t-1}) = \text{Cauchy}(z_t | \sqrt{1-\beta_t} z_{t-1}, \sqrt{\beta_t}).$$

Select  $\beta_t$  s.t.  $\beta_1 < \beta_2 < \dots < \beta_T$ .

Diffusion kernel:

Does it follow that

$$q(z_t | x) = \text{Cauchy}(z_t | \sqrt{\alpha_t} x, (1-\alpha_t) I)$$

$$\text{with } \alpha_t = \prod_{\tau=1}^t (1-\beta_\tau).$$

We proceed via diffusion.

This clearly holds when  $t=1$ , and  $z_0 = x$ .

Assume it holds, that is

$$z_t = \sqrt{\alpha_t} X + (1-\alpha_t) \epsilon_t \quad \text{and} \quad q(z_t | x) = \text{Cauchy}(\sqrt{\alpha_t} x, (1-\alpha_t))$$

$$\epsilon_t \sim \text{Cauchy}(\epsilon_t | 0, 1).$$

Then

$$\begin{aligned} z_{t+1} &= \sqrt{1-\beta_{t+1}} z_t + \sqrt{\beta_{t+1}} \epsilon_{t+1} \\ &= \sqrt{1-\beta_{t+1}} \left( \sqrt{\alpha_t} X + (1-\alpha_t) \epsilon_t \right) + \sqrt{\beta_{t+1}} \epsilon_{t+1} \\ &= \sqrt{\alpha_{t+1}} X + \underbrace{\sqrt{1-\beta_{t+1}} (1-\alpha_t) \epsilon_t}_{\text{Cauchy}(S|0, \sqrt{1-\beta_{t+1}}(1-\alpha_t))} + \underbrace{\sqrt{\beta_{t+1}} \epsilon_{t+1}}_{\text{Cauchy}(S|0, \sqrt{\beta_{t+1}})} \end{aligned}$$

not induction  
sum does not work

It is not true. But we can then define

### Cauchy diffusion

$$z_1 = (1-\beta_1) x + \beta_1 \epsilon_1, \quad \epsilon_1 \sim \text{Cauchy}(\epsilon_1 | 0, 1).$$

$$z_t = (1-\beta_t) z_{t-1} + \beta_t \epsilon_t \quad \epsilon_t \sim \text{Cauchy}(\epsilon_t | 0, 1).$$

Now for the diffusion kernel:

$$t=1, \text{ does } q(z_1 | x) = \text{Cauchy}(z_1 | \alpha_0 x, (1-\alpha_1)).$$

$$\hookrightarrow z_1 = (1-\beta_1) x + \beta_1 \epsilon_1 = \alpha_1 x + (1-\alpha_1) \epsilon_1$$

For  $t \neq 1$ , assuming it.

$$\begin{aligned} z_{t+1} &= (1-\beta_{t+1}) z_t + \beta_{t+1} \epsilon_{t+1} \\ &= (1-\beta_{t+1}) (\alpha_t x + (1-\alpha_t) \epsilon_t) + \beta_{t+1} \epsilon_{t+1} \\ &= \alpha_{t+1} x + \underbrace{(1-\beta_{t+1})(1-\alpha_t) \epsilon_t}_{**} + \underbrace{\beta_{t+1} \epsilon_{t+1}}_{**} \end{aligned}$$

- \*  $(1 - \beta_{t+1})(1 - \alpha_t) \epsilon_t \sim \text{Cauchy}(s \mid 0, (1 - \beta_{t+1})(1 - \alpha_t))$ .
- \*  $\beta_{t+1} \epsilon_{t+1} \sim \text{Cauchy}(s \mid 0, (\beta_{t+1}))$ .
- \* and  $\epsilon_t$  and  $\beta_{t+1} \epsilon_{t+1}$  are independent, and so by Cauchy properties.

2. If  $X \sim \text{Cauchy}(x_0, y_0)$  and  $Y \sim \text{Cauchy}(x_1, y_1)$  are independent, then

$$X + Y \sim \text{Cauchy}(x_0 + x_1, y_0 + y_1)$$

$$X - Y \sim \text{Cauchy}(x_0 - x_1, y_0 + y_1).$$

Hence,

$$\begin{aligned} \epsilon_t + \beta_{t+1} \epsilon_{t+1} &\sim \text{Cauchy}(s'' \mid 0, (1 - \beta_{t+1})(1 - \alpha_t) + \beta_{t+1}) \\ &= \text{Cauchy}(s'' \mid 0, 1 - (1 - \beta_{t+1})\alpha_t) \\ &= \text{Cauchy}(s'' \mid 0, 1 - \alpha_{t+1}). \end{aligned}$$

and so, reparameterizing appropriately, we have that

$$q(z_t | x) = \text{Cauchy}(z_t \mid \mathbb{E} z_t | x, 1 - \alpha_t).$$

Makes sense.

Now, Reverse Process. What is, if any, the distribution such that approximates  $q(z_{t+1} | z_t)$  for small  $\beta_0$ ?

$$q(z_{t-1} | z_t) \propto q(z_t | z_{t-1}) q(z_{t-1}).$$

$$= \frac{1}{\pi \beta_t \left[ 1 + \left( \frac{z_t - (1 - \beta_t) z_{t-1}}{\beta_t} \right)^2 \right]} \cdot q(z_{t-1}).$$

Taylor

$$q(z_{t-1}) \approx q_{z_{t-1}}(z_t) + (z_{t-1} - z_t) D q_{z_{t-1}}(z_t) + \left. (z_{t-1} - z_t)^2 \frac{\partial^2 q_{z_{t-1}}}{\partial z_{t-1}^2} \right|_{z_t}$$

$$\approx \frac{q_{z_{t-1}}(z_t) + (z_{t-1} - z_t) D q_{z_{t-1}}(z_t)}{\pi \beta_t \left[ 1 + \left( \frac{z_t - (1 - \beta_t) z_{t-1}}{\beta_t} \right)^2 \right]}$$

Maybe I need  
2nd order appr. instead.

$$\begin{aligned} \log q(z_{t-1} | z_t) &= \log q(z_t | z_{t-1}) + \log q(z_{t-1}) - \log q(z_t) \\ &= -\log \pi \beta_t - \log \left( 1 + \left( \frac{z_t - (1 - \beta_t) z_{t-1}}{\beta_t} \right)^2 \right) + \log q(z_{t-1}) + \text{const.} \end{aligned}$$

$$\approx -\log \beta_t - \log \left( 1 + \left( \frac{z_t - (1 - \beta_t) z_{t-1}}{\beta_t} \right)^2 \right) + \log q(z_t) + (z_{t-1} - z_t) D(z_t) + \text{const.}$$

$$= -\log \left( \beta_t^2 + (z_t - (1 - \beta_t) z_{t-1})^2 \right) + z_{t-1} D(z_t) + \text{const.}$$

$$= -\log \left( \frac{(z_{t-1}) D(z_t)}{\beta_t^2 + (z_t - (1 - \beta_t) z_{t-1})^2} \right)$$