

$$\Lambda^{-1} - \Theta \Delta \Theta^T = \text{diag}((\Sigma_i^{-1} + R_i^{-1})) - \Theta \Delta \Theta^T$$

α -STABLE DIFFUSIONS:

Extreme Diffusion.

#1.

Grant Strom Peter Hart

$$\Theta \Delta \Theta^T = \text{diag}(\langle S \Sigma_i^{-1} + R_i^{-1} \rangle) - \Theta \Delta \Theta^T$$

The main goal of this notebook lies in developing ~~new~~
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 Levy 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,

$$\psi(t; \alpha, \beta, c, \mu) = \exp[i\mu - |ct|^\alpha(1 - i\beta \operatorname{sign}(t))\phi]$$

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

$\alpha \in (0, 2]$: shape ~ ~~asymmetry~~ concentration.
 $\beta t^{1-\alpha}$: shape - asymmetry
 $c t^{(0, \infty)}$: scale
 $\mu \in (-\infty, \infty)$: location

Properties

A - H A G

It can be expressed
as the sum of an
arbitrary # of
i.i.d r.v's.

- All stable distributions are infinitely divisible.
- With the exception of the normal distribution ($\kappa=2$) stable distributions are leptokurtic and heavy-tailed distributions
- Closure under convolution.

\hookrightarrow can be seen from multiplying chars.

functions.

λ -stable distributions generalize the CT!

Generalized CT.

A non-degenerate random variable \tilde{Z} is λ -stable for some $0 < \lambda \leq 2$ 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

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

$$\Theta \Delta \Theta' = \text{diag}((s\sigma_i^{-1} + r_i^{-1})) - \Theta \Delta \Theta^T$$

Cauchy Distribution.

A relatively small number of κ -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 $X \sim \text{Cauchy}(x_0, y)$ then $kX + \ell \sim \text{Cauchy}(x_0 k + \ell, y|k|)$.

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

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

3. If $X \sim \text{Cauchy}(0, y)$ then $\frac{1}{X} \sim \text{Cauchy}(0, \frac{1}{y})$.

4. Has univariate pdf

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

Has multivariate pdf:

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

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

5. Has univariate cdf:

and multivariate cdf $Q_{\mathbf{x}}(t) = \mathbb{E}[e^{i\mathbf{x}^T t}] = e^{i\mathbf{x}_0^T t - \frac{1}{2} \mathbf{x}^T \Sigma t^2}$.
real functions of degree one s.t. $x_0(t) = \alpha x_0(t)$
 $y(t) = |\alpha| y(t)$.

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

$$KL(Q_{x_0, y_0, \Sigma_1} : Q_{x_0, y_0, \Sigma_2}) = \log \frac{(y_1 + y_2)^2 + (x_{0,1} - x_{0,2})^2}{4 y_1 y_2}$$

Motivating Question:

DDPM's algorithm for training is:

1. Repeat
2. $X_0 \sim q(X_0)$ \mapsto sample from data.
3. $t \sim \text{Unif}(\{\xi_2, \dots, T\})$. $\mapsto T$ is # of timesteps.
4. $\epsilon \sim N(0, I)$.

5. Take gradient step on.

$$\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(X_t, t)\|^2 \mapsto \nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$$

6. Until converged

A crucial step, (5), ~~exist~~ is dependent upon the stability of the normed, because for some schedule of β_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~~ have also a few things that worry me:

- ① - Stability might not be enough? Maybe ~~conjugacy~~ is related to stability, right?
too β needed? Maybe ~~conjugacy~~ is related to stability, right?
- ② - Will it be enough to endow μ and σ to X_0 , and λ ?
(what?)
- ③ - Ho's DDPM claims that for β_t small, the reverse process is also Gaussian. He cites (2015) Sohl-Dickstein for the claim, who in turn cites Feller. (an old article)
How does this claim hold for the Cauchy case?

$$\lambda' = \alpha \circ \Theta^T = f_{\text{diag}}((\sum_i x_i P_i)^T) = \Theta \circ \alpha^T$$

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

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

C

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 volume of (x_1, x_2) is approximately Gaussian when β is small. If I can reproduce this, then I might be able to reproduce it for the laundry case.

Forward forward.

$$x_1 = \sqrt{1-\beta_1} X + \sqrt{\beta_1} \epsilon_1, \quad \epsilon_1 \sim N(\epsilon_1, 0, I).$$

then

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

$$y_t \text{ generally } x_t = \sqrt{1 - \beta_t} z_{t-1} + \sqrt{\beta_t} e_t$$

$$q(\lambda_t | \mathcal{Z}_{t-1}) = N(\lambda_t | \sqrt{1-\beta_t} \mathcal{Z}_{t-1}, \beta_t I)$$

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

Anyways, Bishop's discussion:

Diffusion Kernel.

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

One can marginalize over intermediate variables and find that:

$$q(z_{t+1} | X) = \mathcal{N}(z_t | \sqrt{\alpha_t} X, (1 - \alpha_t) I)$$

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

$$\text{and as } T \rightarrow \infty, \text{ we have}$$

$$q(z_T | X) = \mathcal{N}(z_T | 0, I).$$

$$\text{and} \\ q(z_T) = \mathcal{N}(z_T | 0, I)$$

Indep. of \$X\$.

Conditional distribution.

We are trying to somehow 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

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

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

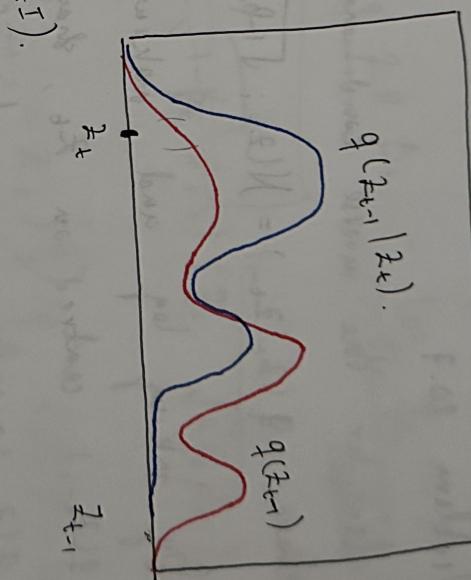
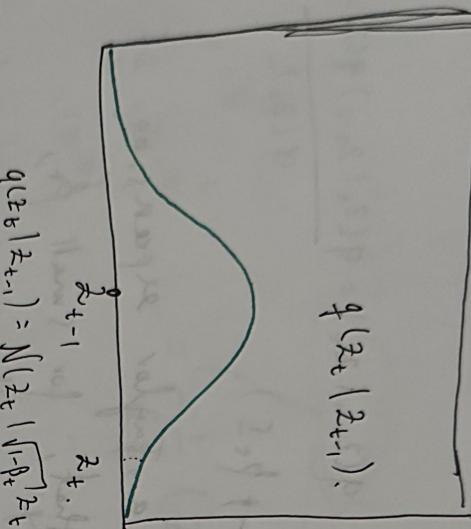
Hence $q(z_{t-1} | z_t)$ is tractable, so we condition on X .

The prior is $\sim q(z_{t-1} | z_t, X) = \frac{q(z_t | z_{t-1}, X) q(z_{t-1} | X)}{q(z_t | X)}$
 left to multiply $q(z_{t-1} | z_t, X) = \frac{(1 - \alpha_{t-1}) \sqrt{1 - \beta_t} z_t + \sqrt{\alpha_{t-1}} \beta_t X}{(1 - \alpha_t) \sqrt{1 - \beta_t} z_t + \sqrt{\alpha_{t-1}} \beta_t X}, \quad \beta_t = \frac{(1 - \alpha_{t-1})}{1 - \alpha_t} I$

$$\Theta \Delta \Theta^T = \text{diag}((s\Sigma_i^{-1} + P_i^{-1})) - \Theta \Delta \Theta^T$$

Reverse Decoder.

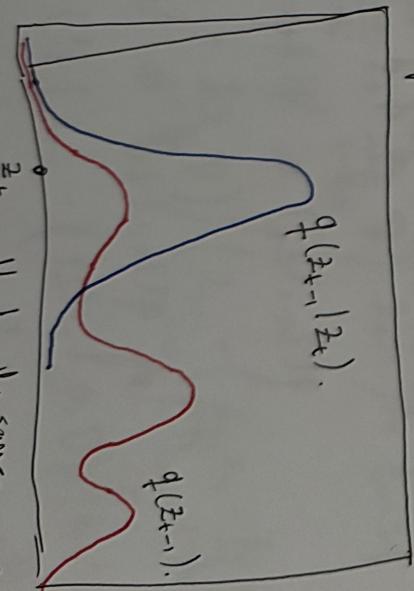
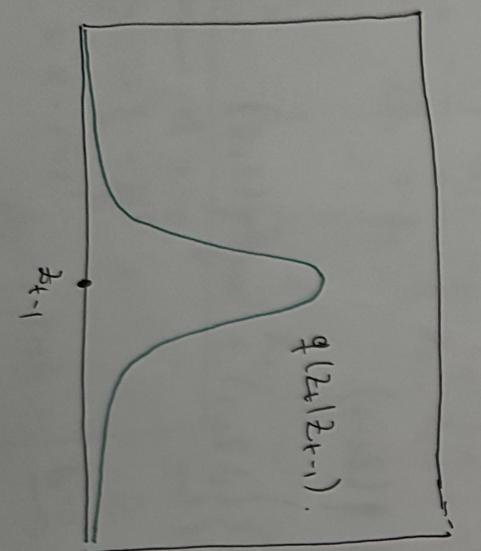
But we really want $q(z_{t-1} | z_t)$ because if we are doing generative modeling beginning from z_T , we can not assume the existence of X . (knowledge)



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

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

If $q(z_t | z_{t-1})$ is sufficiently narrow (Gaussian) then $q(z_{t-1} | z_t)$ 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 works for other works for Causal 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 conditioned

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

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

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

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

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

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

$$\gamma^T = \text{diag}((\Sigma_i^T + P_i^T)) - \Theta \Delta \Theta'$$

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

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

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

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

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

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

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

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

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

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

expansions at $z=0$

$$\begin{aligned} \frac{1}{\sqrt{1-x}} &\approx 1 + \frac{x}{2} + \frac{3x^2}{8} \\ \frac{x}{1-x} &\approx x + x^2 + x^3 + \dots \\ \frac{1}{1-x} &\approx 1 + x + x^2 + \dots \end{aligned}$$

Can we apply this to a Cauchy diffusion?

Let's define A :

Forward Finsler:

Consider X given, in the univariate setting.

$$Z_1 = \sqrt{1-\beta_1} X + \sqrt{\beta_1} \epsilon_1, \quad \epsilon_1 \sim \text{Cauchy}(\epsilon_1 | 0, 1). \quad (1)$$

$$\Rightarrow q(Z_1 | X) = \text{Cauchy}\left(Z_1 | \sqrt{1-\beta_1} X, \sqrt{\beta_1} I\right)?$$

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

| 1. If \forall_n Cauchy(X_0, Y) then $kX + l \sim \text{Cauchy}(X_0k + l, Y|k|)$.

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

$$Z_1 | X \sim \text{Cauchy}\left(\sqrt{1-\beta_1} X, \sqrt{\beta_1}\right).$$

Each successive margin is given by:

(3)

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

Where again, by the Markov property and 1.

$$q(Z_t | Z_{t-1}) = \text{Cauchy}\left(Z_t | \sqrt{1-\beta_t} Z_{t-1}, \sqrt{\beta_t}\right).$$

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

Diffusion kernel:

Does it follow that

$$q(Z_t | X) = \text{Cauchy}\left(Z_t | \sqrt{\kappa_t} X, (1-\kappa_t) I\right)$$

$$\text{with } \kappa_t = \prod_{r=1}^t (1-\beta_r).$$

$$-\alpha_i \omega_i^T = \text{diag}((\zeta \Sigma_i^{-1} + R_i^{-1})) - \Theta \Delta \Theta^T$$

We proceed via diffusion.

This clearly holds when $t=2$, and $x_0 = x$.

Assume it holds, that is
 $x_t = \sqrt{\alpha_t} X + (1-\alpha_t) e_t$ and $q(x_t | X) = \text{Cauchy}(\sqrt{\alpha_t} X, (1-\alpha_t))$

$$x_{t+1} = \sqrt{1-\beta_{t+1}} x_t + \sqrt{\beta_{t+1}} e_{t+1}$$

$$\begin{aligned} &= \sqrt{1-\beta_{t+1}} \left(\sqrt{\alpha_t} X + (1-\alpha_t) e_t \right) + \sqrt{\beta_{t+1}} e_{t+1} \\ &= \sqrt{\alpha_{t+1}} X + \underbrace{\sqrt{1-\beta_{t+1}} (1-\alpha_t) e_t}_{\text{Cauchy}(s|0, \sqrt{1-\beta_{t+1}} (1-\alpha_t))} + \sqrt{\beta_{t+1}} e_{t+1} \end{aligned}$$

It is not true. But we can then define

- sum does not work

Cauchy diffusion

$$x_1 = (1-\beta_1) X + \beta_1 e_1, \quad e_1 \sim \text{Cauchy}(e_1 | 0, \underline{\lambda})$$

$$x_t = (1-\beta_t) x_{t-1} + \beta_t e_t \quad e_t \sim \text{Cauchy}(e_t | 0, \underline{\lambda}).$$

$$x_{t+1} \sim \text{Cauchy}(x_t | (1-\beta_t)x_{t-1}, \beta_t)$$

Now for the diffusion kernel:

$$t=1, \quad \text{does } q(x_1 | X) = \text{Cauchy}(x_1 | \alpha_1 X, (1-\alpha_1))$$

$$\hookrightarrow x_1 = (1-\beta_1) X + \beta_1 e_1 = \alpha_1 X + (1-\alpha_1) e_1$$

For $t+1$, assuming t :

$$\begin{aligned} x_{t+1} &= (1-\beta_{t+1}) x_t + \beta_{t+1} e_{t+1} \\ &= (1-\beta_{t+1}) (\alpha_t X + (1-\alpha_t) e_t) + \beta_{t+1} e_{t+1} \\ &= \alpha_{t+1} X + (1-\alpha_{t+1})(1-\alpha_t) e_t + \beta_{t+1} e_{t+1} \end{aligned}$$

$$(1 - \beta_{t+1})(1 - \alpha_t) \epsilon_t \sim \text{Cauchy}(s | 0, (1 - \beta_{t+1})(1 - \alpha_t)).$$

$\beta_{t+1}, \alpha_{t+1} \sim \text{Cauchy}(s'' | 0, (\beta_{t+1}))$.
and ϵ_t 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

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

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

Hence,

$$x + \alpha_t \sim \text{Cauchy}(s'' | 0, (1 - \beta_{t+1})(1 - \alpha_t) + \beta_{t+1})$$

$$= \text{Cauchy}(s'' | 0, 1 - (1 - \beta_{t+1})\alpha_t)$$

$$= \text{Cauchy}(s'' | 0, 1 - \alpha_{t+1}).$$

and so, reparameterizing appropriately, we have that

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

Moreso.

Now, Densoge process. What is, if any, the distribution such that approximates $q(z_{t+1} | z_t)$ for small β_{t+1} ?

$$= S \cdot \mathcal{H}^n \cdot \left(\sum_i k_i \right) \otimes I_P$$

$$g(z_{t+1} | z_t) \propto g(z_t | z_{t+1}) \cdot q(z_{t+1}).$$

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

Taylor

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

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

Number \leq need O
2nd order approx.

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