

# Diffusion Models (ch. 20 Bishop)

Wednesday, November 13, 2024 7:12 AM

## §20 Diffusion Models (still following Bishop's book)

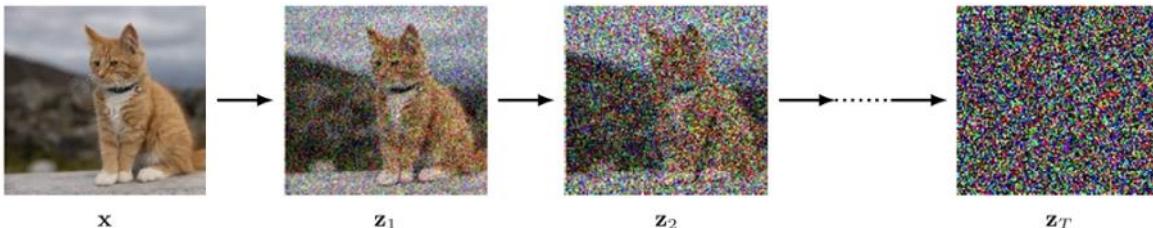
aka "Denoising diffusion probabilistic models" or DDPM

Sohl-Dickstein et al. '15  
Ho, Jaini, Abbeel '20

Other generative approaches: GANs, VAE, Normalizing Flow

Diffusion models: easier to train, but slow to generate new samples  
State-of-the-art!

Inspiration: we corrupt training data



$$z_t = \sqrt{1-\beta_t} x + \sqrt{\beta_t} \varepsilon_t \xrightarrow{\varepsilon_t \sim N(0, I)} \text{pushes mean toward } 0, \text{ variance toward } I$$

$$z_t = \sqrt{1-\beta_t} z_{t-1} + \sqrt{\beta_t} \varepsilon_t \quad \begin{matrix} \text{In lit.,} \\ "x_t" \text{ is our "x"} \\ "z_t" \text{ is our "z_t"} \end{matrix} \quad 0 < \beta_1 < \beta_2 < \dots < \beta_T < 1 \quad \begin{matrix} \text{chosen by hand} \end{matrix}$$

$$\boxed{z_{t-1}} \xrightarrow{g(z_t | z_{t-1})} \boxed{z_t} \quad \begin{matrix} \text{(Fixed)} \\ g(z_t | z_{t-1}) = N(\sqrt{1-\beta_t} z_{t-1}, \beta_t I) \\ \text{"Forward Process" (like "encoder" in VAE,} \\ \text{but not learned)} \end{matrix}$$

Altogether,  $(z_t)$  is a Markov Chain. Conditioned on starting point (observed data  $x$ ),

$$g(z_1, \dots, z_T | x) = g(z_1 | x) \cdot \prod_{t=2}^T g(z_t | z_{t-1})$$

so ("marginalize" over  $z_1, \dots, z_{T-1}$ )

$$g(z_T | x) \sim N(\sqrt{\alpha_T} x, (1-\alpha_T) I)$$

$$\alpha_T = \prod_{t=1}^T (1-\beta_t) \rightarrow 0 \quad \begin{matrix} \text{as } T \rightarrow \infty \\ (\text{since } 0 < \beta_t < 1) \end{matrix}$$

"diffusion kernel" (can do this for any  $T$ . Useful for training as we can jump to any  $T$  w/o starting at beginning)

Goal: learn to undo this forward process

Then we can sample  $z_\infty$  (i.e.  $z_T$  for  $T$  large) from  $N(0, I)$ , go backward, and end w/ a sample  $x \sim p(x)$

$$g(z_{t-1} | z_t) = \frac{g(z_t | z_{t-1})}{g(z_t)} \frac{g(z_{t-1})}{g(z_t)} \quad \begin{matrix} \text{known} \\ \text{intractable} \end{matrix} \quad g(z_{t-1}) = \int g(z_{t-1} | x) p(x) dx \quad \begin{matrix} \text{unknown} \end{matrix}$$

so instead look at

$$g(z_{t-1} | z_t, x) = \frac{g(z_t | z_{t-1})}{g(z_t)} \cdot g(z_{t-1} | x) \quad \begin{matrix} \text{irrelevant due to Markov Structure} \\ \text{diffusion kernel} \end{matrix}$$

$\frac{g(z_t | z_{t-1})}{g(z_t)}$  ← irrelevant for now

$$\text{For training, } |x| \text{ is ok} \quad = N(m_t, \sigma_t^2 I) \quad m_t = \frac{((1-\alpha_{t-1})\sqrt{1-\beta_t} z_t + \sqrt{1-\alpha_{t-1}} \beta_t x_t)}{1-\alpha_t}$$

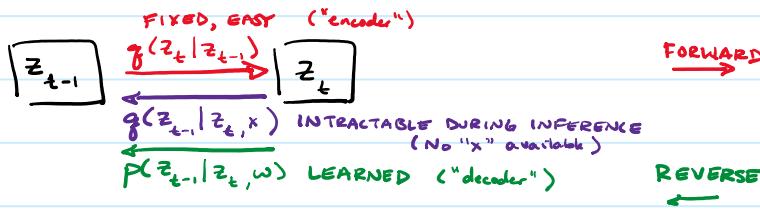
For inference, it isn't!

$$\sigma_t^2 = \frac{\beta_t (1-\alpha_{t-1})}{1-\alpha_t}$$

# Diffusion Models (p. 2)

Wednesday, November 13, 2024 9:52 AM

Reverse Decoder we'll learn model parameters  $\omega$  to approximate  $q(z_{t-1}|z_t)$  the reverse distribution.



Keep  $\beta_t \ll 1$  so forward encoding is small (almost identity) so easier to invert.

Then use many steps ( $T \approx 1000's$ ) so  $\alpha_t \approx 0$

We'll parameterize

$$p(z_{t-1}|z_t, \omega) = N(\mu, \beta_t I)$$

$\mu(z_t, \omega, t)$  via neural net.

often (not always) fixed.

Overall,

$$p(x, z_1, \dots, z_T | \omega) = p(z_T) \cdot \underbrace{\prod_{t=2}^T p(z_{t-1}|z_t, \omega)}_{\substack{\sim N(0, I) \\ \text{seed}}} \cdot \underbrace{p(x|z_1, \omega)}_{\substack{\text{neural net eval}}}$$

Conceptually straightforward

(but  $T \approx 1000's$  makes it slow)

## Training

$\int p(\vec{z})$  intractable

Objective is likelihood,  $p(x|\omega) = \int \dots \int p(x, z_1, \dots, z_T | \omega) dz_1 \dots dz_T$  via neural net Intractable

Overcome/approximate w/ standard trick

("the" trick of "Variational Inference", used in VAE and EM) see also fig. 6.3

## "Evidence Lower BOund" (ELBO)

Def: The Kullback-Leibler divergence of two pdfs  $f, g$

$$KL(f \| g) = - \int f(z) \cdot \ln \left( \frac{g(z)}{f(z)} \right) \quad \text{see § 2.5.5}$$

see § 2.5.7

⚠ Not symmetric in  $f, g$

Thm  $KL(f \| g) \geq 0$  (and  $= 0$  iff  $f = g$  a.e.)

proof sketch: Use convexity of  $-\log(\cdot)$  w/ Jensen's inequality.

Then

$$\ln(p(x|\omega)) = \underbrace{\int g(\vec{z}) \cdot \ln \left( \frac{p(x, \vec{z}|\omega)}{g(\vec{z})} \right) d\vec{z}}_{L(\omega)} - \underbrace{\int g(\vec{z}) \cdot \ln \left( \frac{p(\vec{z}|x, \omega)}{g(\vec{z})} \right) d\vec{z}}_{\text{IGNORE}}$$

"evidence lower bound"

tractable since we know  $g$   
(see Eq. 20.26, use M.C.)

$$= KL(g, p(\vec{z}|x, \omega))$$

$\geq 0$  if

$$g(\vec{z}) = p(\vec{z}|x, \omega)$$

want to maximize this...

" $\vec{z}$ " really means  $z_1, \dots, z_T$

$\int g(\vec{z})$  tractable

.. instead, maximize this (lower bound).

(we choose  $g$  to make tractable, but also attempt ↗)

# Diffusion Models (p. 3)

Wednesday, November 13, 2024 10:36 AM

(improvement: predicting noise, § 20.2.4) Ho, Jain, Albeeel '20

... plus simplifications, approximations...

Ultimately,  $\text{loss}(\omega) = -\sum_{t=1}^T \|\underbrace{g(\sqrt{\alpha_t}x + \sqrt{1-\alpha_t}\varepsilon_t, \omega, t)}_{\text{predicted noise}} - \varepsilon_t\|_2^2$

Inputs  
predicted noise  
neural net  
training data  
weights (per training sample)  
in practice, subsample (SGD style) during training

Inference: generating new samples  $x \sim p(x)$

$$z_T \sim N(0, I)$$

for  $t = T, T-1, \dots, 2$

$$\alpha_t = \prod_{i=1}^t (1 - \beta_i)$$

$$\mu(z_t, \omega, t) \leftarrow \frac{1}{\sqrt{1-\beta_t}} \left( z_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} g(z_t, \omega, t) \right)$$

$$\varepsilon \sim N(0, I)$$

$$z_{t-1} \leftarrow \mu(z_t, \omega, t) + \sqrt{\beta_t} \cdot \varepsilon$$

$$x = \frac{1}{\sqrt{1-\beta_1}} \left( z_1 - \frac{\beta_1}{\sqrt{1-\alpha_1}} g(z_1, \omega, t) \right)$$

slow since  $T \approx 1000$ ;

Improvement: view as solving stochastic diff. eq.  
and use better discretization

Relation to Score-matching (Hyvärinen '05, Song, Erman '19)

Score score  $S(x) = \nabla_x \ln p(x)$  (ignores multiplicative normalization)

$$\text{loss } J(\omega) = \frac{1}{2} \int \|\underbrace{S(x, \omega) - \nabla_x \ln p(x)}_{\text{approximate score fun via neural net}}\|_2^2 p(x) dx$$