# 4/28 - Probabilistic physics w/ variational inference

- Last day of class → wed 4/28
    - → everyone should be making progress
    - → last chance for feedback
    - → final presentations 5/9 3-5 in DRLB-A2
    - → report due by 5/12

- Final Course Feedback
    - Please be honest - first time class and feedback will drastically reshape next iteration
        - → More HW
        - → Stricter pre-regs

---

Recall ELBO from last time

$$\mathcal{E} = -\mathbb{E}_{x \sim z}\left[\log p(x|z)\right] - KL\left(g(z|x) \| p(z)\right)$$

We discussed Kingma + Welling "vanilla" VAE

single-sample Monte Carlo

① $\mathbb{E}_{x \sim z}\left[\log p(x|z)\right] \approx \log p(x_d | z_d)$

② $p(z) = N\left(\mu=0, \Sigma = I\right)$

Which allowed generative modeling

① Sample $z \sim N(0, I)$

② Decode using $p(x | z)$

There is a craft to designing architectures around different choices of embeddings $z$, and priors

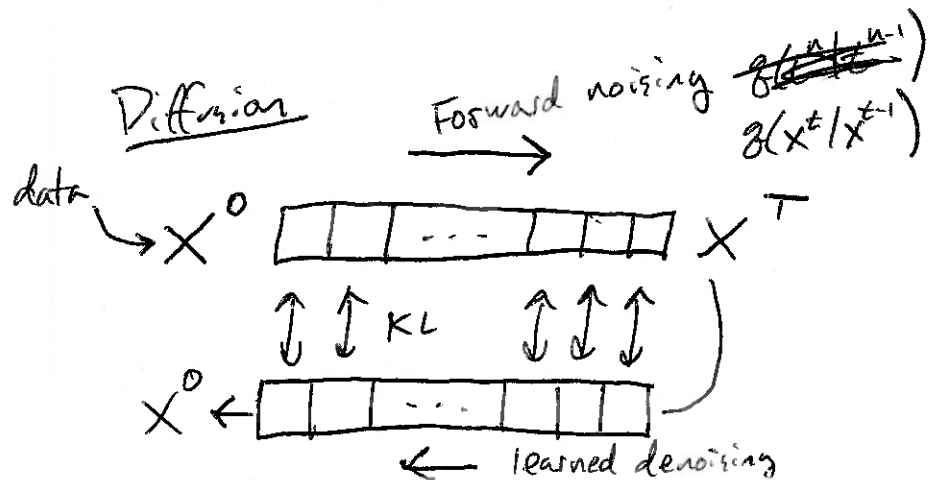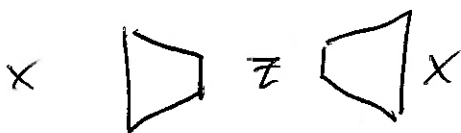Two examples

- Denoising Diffusion
- Physical Priors

GOAL
Computationally tractable
ELBO!

↑
Gaussian
Gaussian
Gaussian!

Denoising Diffusion          (Sohl-Dickstein)

___

Idea    Don't jump from $x$ to $z$, instead decode in increments

Vanilla

$x$ ▷ $z$ ◁ $x$

Diffusion

Forward noising $\cancel{q(x^{n} | x^{n-1})}$
$q(x^t | x^{t-1})$

data $\to x^0$ [ | | ⋯ | | ] $x^T$

↕ ↕ KL   ↕ ↕ ↕

$x^0 \leftarrow$ [ | | ⋯ | | ]

$\leftarrow$ learned denoising

# Ingredients

Forward/Noising

$q(x_0)$ — data distribution

$$q(x^t | x^{t-1}) = T(x^t | x^{t-1}; \beta_t) \qquad \text{a kernel to add noise}$$

we'll a-tsek to simple →
$$= N(x^t; x^{t-1}\sqrt{1-\beta_t}, \beta_t I)$$

which is equiv to
$$x_t = \sqrt{1-\beta_t}\, x_{t-1} + \sqrt{\beta_t}\, \varepsilon \qquad \leftarrow N(0,I)$$

And we get joint dist
$$q(x^0, \dots x^T) = q(x^0) \prod_{t=1}^{T} q(x^t | x^{t-1}) \qquad \leftarrow \text{gaussian increments}$$

Lemma
$$\lim_{T \to \infty} x_t = N(0, I)$$

Note   Before we prove this lemma, note that we get a unit gaussian embedding __by construction__ instead of using KL to get it __by__ __penalty__

**Pf** To prove by induction, check 2 steps

$$X_t = \sqrt{1-\beta_t} \, X_{t-1} + \sqrt{\beta_t} \, \varepsilon_t$$

$$X_{t+1} = \sqrt{1-\beta_{t+1}} \sqrt{1-\beta_t} \, X_{t-1} + \underbrace{\sqrt{\beta_{t+1} \beta_t} \, \varepsilon_t + \sqrt{\beta_{t+1}} \, \varepsilon_{t+1}}$$

$$\underbrace{\phantom{\sqrt{1-\beta_{t+1}} \sqrt{1-\beta_t} X_{t-1}}}$$

products          Additive gaussian noise
w/ weights

Let $\alpha_s = 1 - \beta_s$

$$\overline{\alpha}_t = \prod_{s=1}^{t} \alpha_s$$

By induction can show

$$X_t = \sqrt{\overline{\alpha}_t} \, X_0 + \sum_{s=1}^{t} \left\{ \sqrt{\overline{\alpha}_{s-1} \beta_s} \, \varepsilon_s \right.$$

For gaussians, variance of sum is sum of var

$$X_t \sim N\left( \sqrt{\overline{\alpha}_t} \, X_0 \,, \, \sum_{s=1}^{t} \overline{\alpha}_{s-1} \beta_s \, I \right)$$

**Lemma** $\quad \sum_{s=1}^{t} \overline{\alpha}_{s-1} \beta_s = 1 - \overline{\alpha}_t$

**Pf** $\quad \overline{\alpha}_s = \overline{\alpha}_{s-1} \alpha_s \qquad\qquad\qquad \overline{\alpha}_{s-1} \beta_s = \overline{\alpha}_{s-1} - \overline{\alpha}_s$

$$= \overline{\alpha}_{s-1} (1 - \beta_s) \quad \Rightarrow$$

So that

$$\sum_s \overline{\alpha}_{s-1} \beta_s = \sum_s \overline{\alpha}_{s-1} - \overline{\alpha}_s = 1 - \overline{\alpha}_t$$

↖ telescoping series

Finally, take limits

Let $0 < \beta_s \ll 1$

$$\lim_{t \to \infty} |\bar{\alpha}_t| = \lim_{t \to \infty} \left| \prod_{s=1}^{t} (1 - \beta_s) \right|$$

$$\leq \lim_{t} \left| \prod_{s=1}^{t} (1 - \beta_{min}) \right|$$

$$\leq \lim_{t} (1 - \beta_{min})^t$$

$$= 0$$

So

$$\lim_{t \to \infty} N\left( \sqrt{\bar{\alpha}_t} \, x_0, \; \sum_{s=1}^{t} \bar{\alpha}_{s-1} \beta_s I \right)$$

$$= \lim_{t \to \infty} N\left( \sqrt{\bar{\alpha}_t} \, x_0, \; 1 - \bar{\alpha}_t \right)$$

$$= N(0, I)$$

Back to forward ingredient

$$q(x_1^0, \dots x^T) = q(x^0) \prod_{t=1}^{T} q(x^t | x^{t-1})$$

$\hookrightarrow$ gaussian!

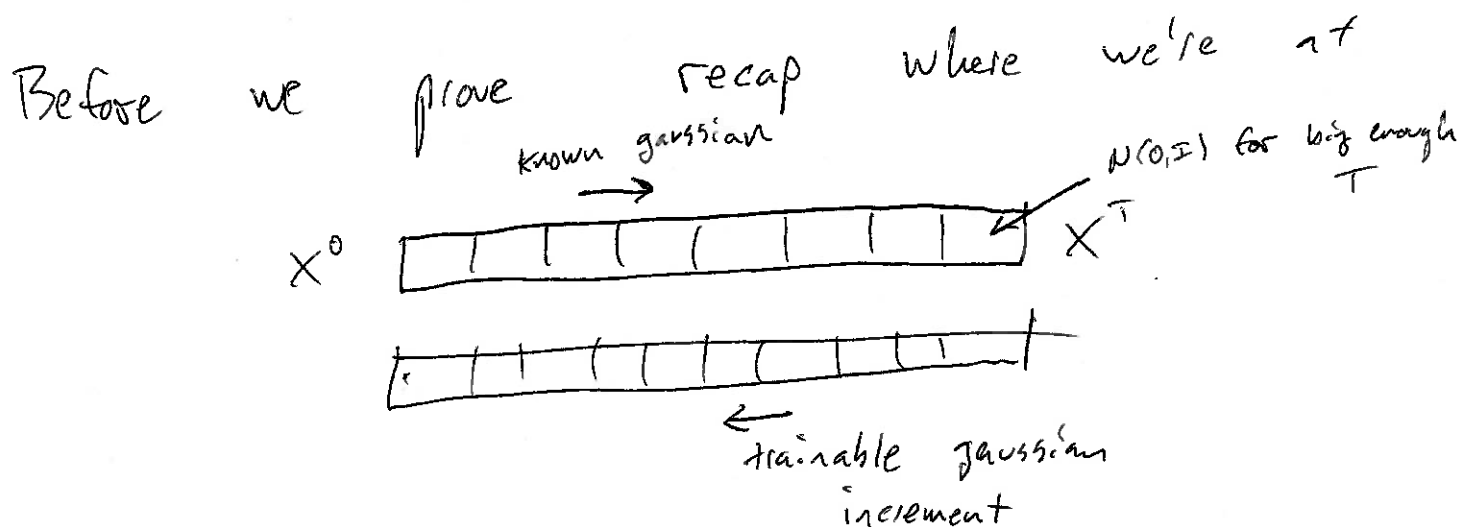# Ingredient 2

$$P(x^T) = \text{prior distribution}$$

$$P(x^0, \ldots, x^T) = P(x^T) \prod_{t=1}^{T} P(x^{t-1} | x^t)$$

going **backward** in time now

<u>Claim</u>  For small $\beta$, $P(x^{t-1} | x^t)$ is Gaussian

we will prove this as justification for
the parameterization

$$P(x^{t-1} | x^t) = N\left( \mu = f(x^t; \theta_f), \ \Sigma = g(x^t; \theta_g) \right)$$

Learnable NN's

Before we prove recap where we're at

known gaussian →

$N(0, I)$ for big enough $T$

$x^0$ | | | | | | | | | $x^T$

| | | | | | | | | |

← trainable gaussian increment

Final step will be to make steps match w/ KL

We'll finally show

$$\cancel{\frac{q(x_{t-1}|x_t, x_0)}{q(} = q}$$

$$q(x_{t-1}|x_t) = q(x_{t-1}|x_t, x_0) + O(\beta_t)$$

## Justification of reverse transition being Gaussian

By Bayes
$$q(x^{t-1} | x^t, x^0) = \frac{q(x^t | x^{t-1}, x^0) \, q(x^{t-1} | x^0)}{q(x^t | x^0)}$$

From forward
$$q(x^t | x^{t-1}) = N\left(\sqrt{1-\beta_t}\, x_{t-1}, \, \beta_t I\right)$$

As we derived already
$$q(x^{t-1} | x^0) = N\left(\sqrt{\bar{\alpha}_{t-1}}\, x_0, \, (1-\bar{\alpha}_{t-1}) I\right)$$

product of Gaussians is an un-normalized Gaussian
for some $M^*, \Sigma^*, C$ that we **could** derive
$$\Rightarrow \quad q(x^t | x^{t-1}) \, q(x^{t-1} | x^0) \sim C \, N(M^*, \Sigma^A)$$

Skipping details
$$q(x_{t-1} | x_t, x_0) = N\left(M_t^*(x_0, x_t), \, \tilde{\beta}_t^* I\right)$$

$$M^* = \frac{\sqrt{\bar{\alpha}_{t-1}}\, \beta_t}{1 - \bar{\alpha}_t} x_0 + \frac{\sqrt{\bar{\alpha}_t}\, (1-\bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t$$

$$\beta^* = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}$$

Q- How do ~~they~~ depend on $\beta_t$ when approximating $q(x_{t-1} | x_t)$?

To remove the conditioning on $X_0$ we'll use the laws of total expectation + variance:

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|Y]]$$

$$\text{var}[X] = \mathbb{E}[\text{var}[X|Y]] + \text{var}[\mathbb{E}[X|Y]]$$

On the mean

$$\mathbb{E}[x_{t-1}|x_t] = \mathbb{E}[\mathbb{E}[x_{t-1}|x_t, x_0]]$$

$$= \mathbb{E}[\mu^*] \sim O(\beta_t^2)$$

On the variance

$$\text{var}[x_{t-1}|x_t] = \mathbb{E}[\text{var}(x_{t-1}|x_t, x_0)] + \text{var}[\mathbb{E}[x_{t-1}|x_t, x_0]]$$

$$= \underbrace{\mathbb{E}[\beta^* I]}_{O(\beta_t^2)} + \underbrace{\text{var}[\mu^*]}_{O(\frac{\beta}{t})}$$

(we're skipping the estimates on $\beta_t$)

~~To summarize~~

# Final Ingredient #3 - the ELBO

$$\varepsilon = \mathbb{E}_{g(x^T | x^0)} \left[ \log p(x_0 | x_T) \right] - \sum_{t=2}^{T} \mathbb{E}_{g(x_t | x_0)} \left[ KL\left( g(x_{t-1} | x_t, x_0) \| p(x_{t-1} | x_t) \right) \right]$$

- Learn to denoise single diffusion steps
- Generate by sampling a gaussian + denoising

## Some takeaways

- Designed an encoding which
    - required no learning
    - took us to a latent gaussian
    - has gaussian increments
- Designed a decoding which
    - is indirectly supervised by encoding
    - gives gaussian increments

- When combined in ELBO
    - Lots of closed form expressions for KL's

How can you use this in your research?
(How much math do you need to track?)

**V1** Off-the-shelf
- just grab one of the implementations
- understand that if you mess w/ the architecture, you may break the theory
  - maybe leave ELBO loss alone
  - don't play w/ input/output of denoising,
  - do play w/ architectures, $\beta$ schedules

**V2** Improve
- Many other physical models can give $\lim_{T \to \infty} X^T \sim N(0, I)$
- Replace noising w/ something physical
- Think about challenges & requirements
  - Expensive to generate - how can we use shorter T
  - What other SDEs allow gaussian reverse process?
  - SPDES?