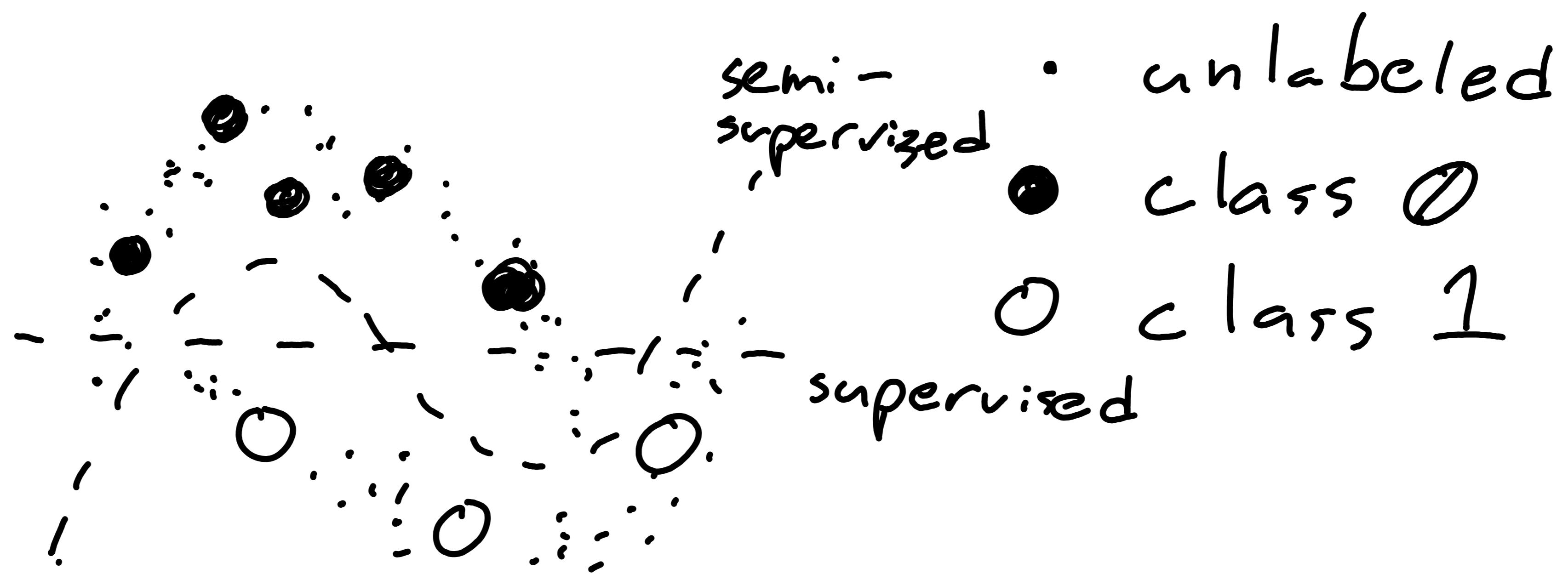


Semi-Supervised Learning

Supervised learning : $x, y \sim p(x, y)$

Semi-supervised

: $x \sim p(x)$ also
unlabeled examples



Simplest approach: (self-training, pseudo-labeling)
Train the model to predict its predictions:

$$L = -\operatorname{argmax}_y [p_0(y|x)] \log p_0(y|x)$$

Problem: Model can reinforce (get more
confident about) its own predictions.

Consistency regularization:

Model's output should be the same when the input is perturbed.

$$L = - \sum_y p_\theta(y|x) \log p_\theta(y|x')$$

original x perturbed x

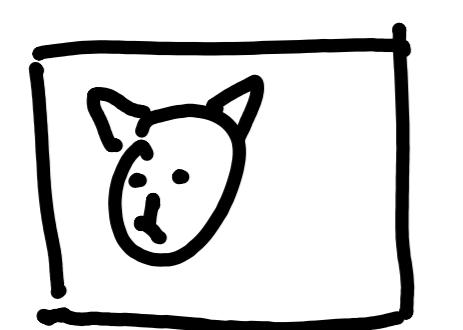
How to perturb x ?

Simple: add a little noise.

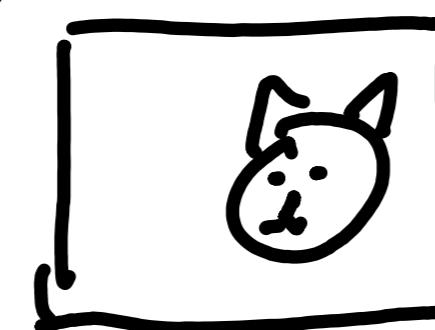
Better: use data augmentation.

→ modify x so that it's still "realistic" but the class is retained.

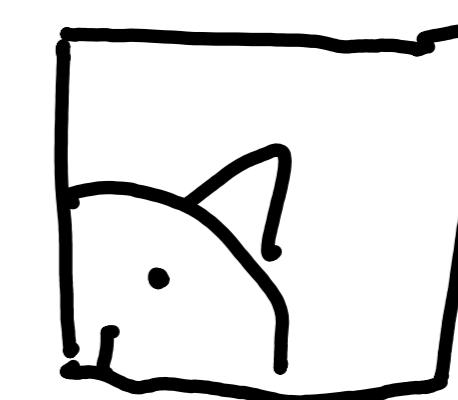
Allows us to bake in domain knowledge.



horizontal
flip



rescale
and crop



manipulate
color



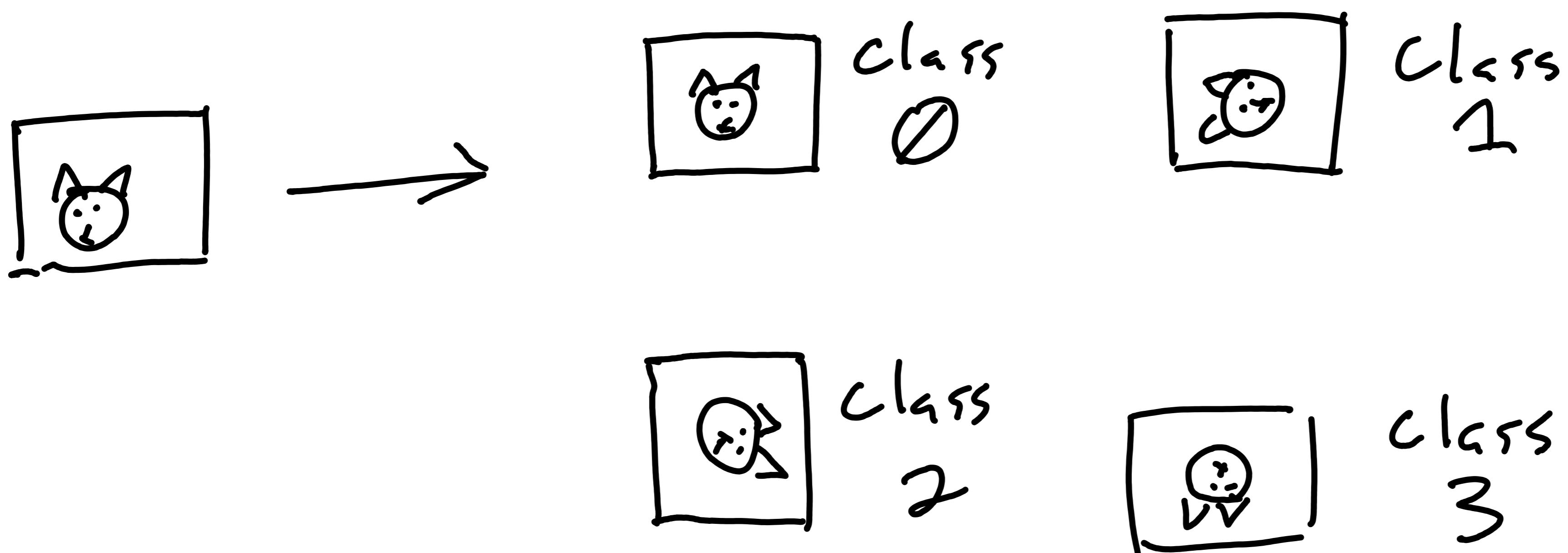
Self-supervised learning

We have $x \sim p(x)$ and can hand-design a function that gives a_y for a given x .

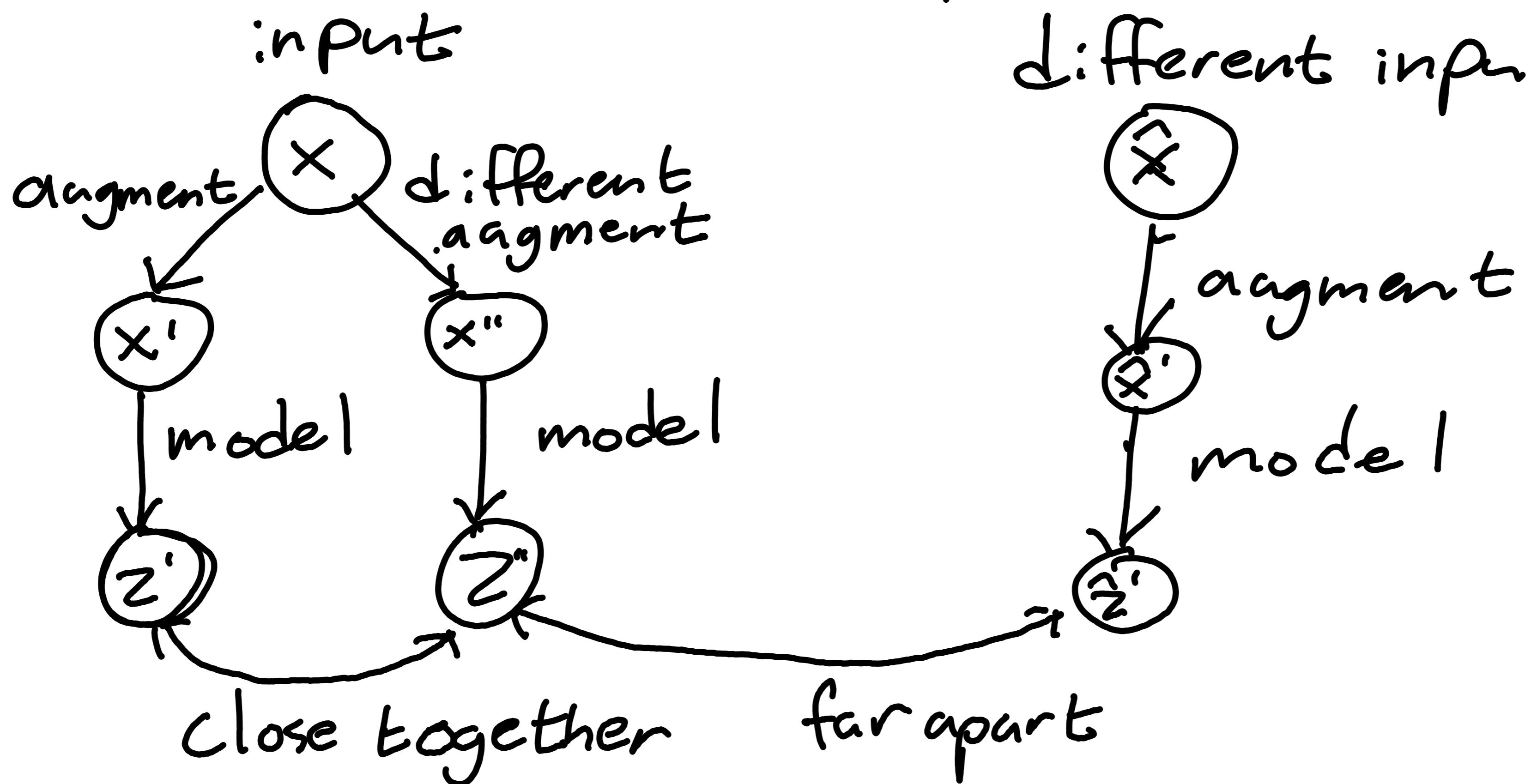
~~ex~~ Masked language modeling

sentence \rightarrow masked sentence \rightarrow sentence
automatic need
no labelling

~~ex~~ Predicting rotations



Contrastive self-supervised learning: (SimCLR)



$$L = -\log \frac{\exp(\text{sim}(z', z'')/\tau)}{\sum_{b=1}^B \exp(\text{sim}(z', \hat{z}_b)/\tau)}$$

hyperparameter
e.g. cosine distance

Generative Modeling

Generative model: Model $p(x)$ with $q_\theta(x)$ based on
 $x \sim p(x)$

- Why?
1. If $q_\theta(x)$ approximates $p(x)$ well,
we can sample $x \sim q_\theta(x)$ and x
will be "realistic".
 2. We can use $q_\theta(x)$ to approximate
the probability of x .
 3. Learning $q_\theta(x)$ can "help" model
 $p(y, x)$ (transfer learning).

Generative Adversarial Networks:

$$\min_{\theta} \max_{\phi} \mathbb{E}_{\substack{x \sim p(x) \\ \hat{x} \sim q_{\theta}(x)}} V(f_{\phi}(x), f_{\phi}^{\top}(\hat{x}))$$

loss function
neural network "discriminator"
neural net "generator"

Advantages:

- f_{ϕ} is a neural network, so we can give our divergence an inductive bias.
- Change V to change behavior.
- You can use anything you want for $q_{\theta}(x)$ as long as you can $\hat{x} \sim q_{\theta}(x)$ and it's differentiable.

Disadvantages:

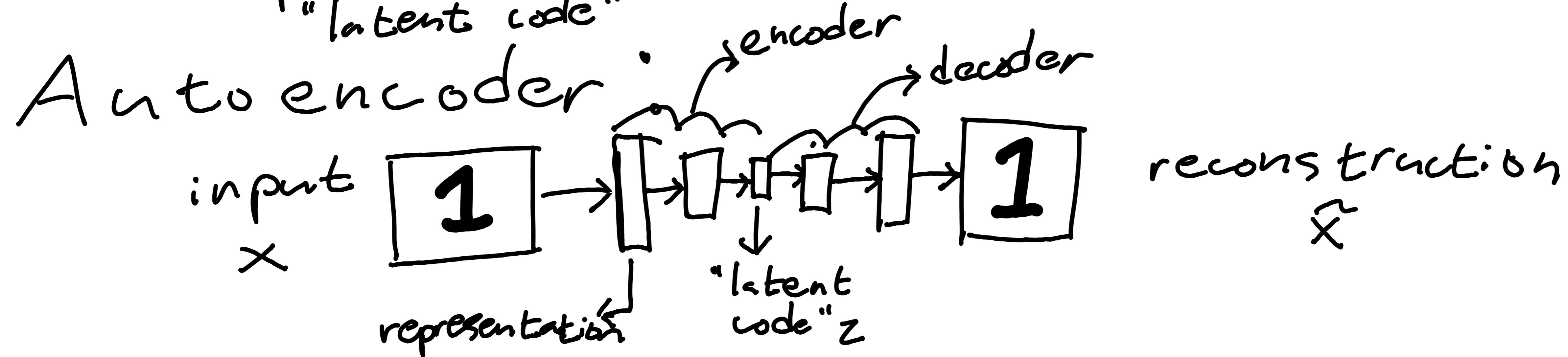
- Train with alternating gradients, unstable.
- Can't actually evaluate $q_{\theta}(x)$

Variational Autoencoder

Generative process:

$$z \sim p(z)$$

prior on
"latent code"



Train using a reconstruction loss

$$\text{e.g. } \|x - \tilde{x}\|_2^2$$

If $\dim(z) < \dim(x)$, we are (maybe) doing compression.

Problem: z is continuous \rightarrow it can store infinite information.

VAE: Regularize the information in z .

Objective = reconstruction + regularization

$$= \mathbb{E}_{z \sim q_{\phi}(z|x)} \log P_{\theta}(x|z) - KL(q_{\phi}(z|x) || p(z))$$

z is our latent code (stochastic).

$q_{\phi}(z|x)$ is "encoder", maps input \rightarrow latent

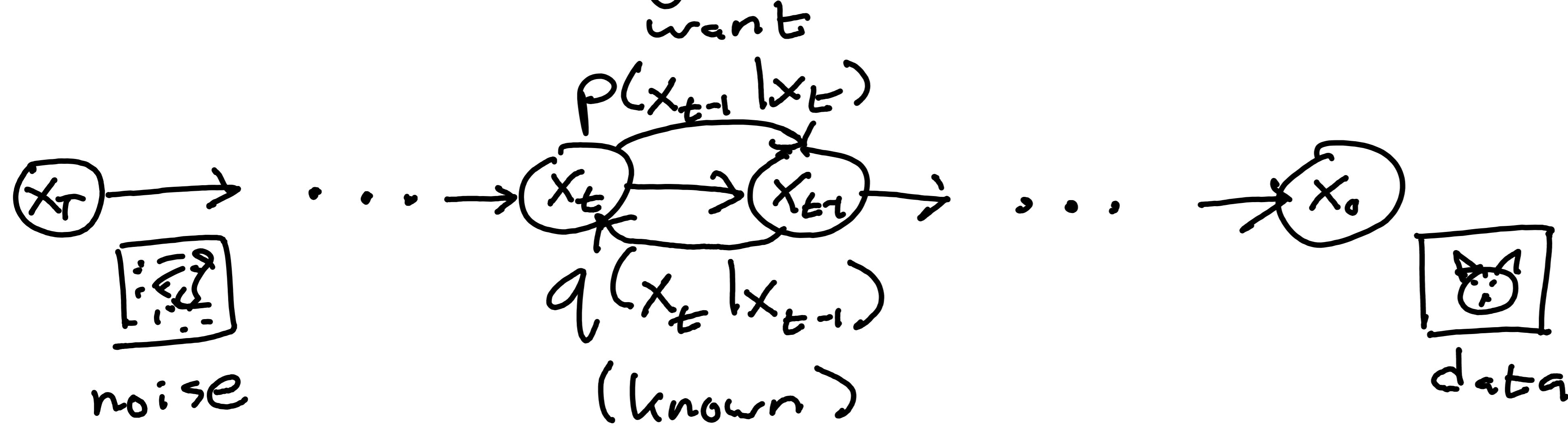
$P_{\theta}(x|z)$ is "decoder", maps latent \rightarrow reconstruction

$p(z)$ is prior over latents.

KL term measures information in z .

Diffusion models

Given $x \sim p(x)$, define a "forward" diffusion process that goes from x to noise.



$$q(x_t | x_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

\downarrow noise schedule
 β_t small

$$q(x_t | x_0) = \mathcal{N}(\bar{\alpha}_t^t x_0, (1 - \bar{\alpha}_t) I)$$

$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i, \quad \alpha_i = 1 - \beta_i$$

Learn a neural network $P_\theta(x_{t-1} | x_t)$

$$P(x_{t-1} | x_t) = \mathcal{N}(\mu_\theta(x_t, t), \beta_t I)$$

↳ optimal under certain assumptions

(If β_t is small enough, $q(x_{t-1} | x_t)$ will be Gaussian)

Now we just need to train μ_θ to predict

$$\hat{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \varepsilon_t \right)$$

↳ noise applied at step t
known during training

$$L_t = \mathbb{E}_{t \sim [1, T]} \| \varepsilon_t - \varepsilon_\theta(x_t, t) \|_2^2$$

$$= \mathbb{E}_{t \sim [1, T]} \| \varepsilon_t - \varepsilon_\theta(\sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \varepsilon_{t-1}, t) \|_2^2$$