

Classifier

Free

guidance

MIT Diffusion

Models

structure

Classifier-Free Guidance

- goal: generate stuff with prompts.

How do we make models fit to prompts.

* Vanilla guidance → vector sets $\begin{cases} \text{N.N} \\ u_t^\theta: \mathbb{R}^d \times \mathbb{Y} \times [0,1] \\ \rightarrow \mathbb{R}^d \\ (x, y, t) \rightarrow u_t^\theta(x|y) \\ x_0 \sim p_{\text{init}} \\ dx_t = u_t^\theta(x_t|y) dt + \sigma dw_t \\ x_1 \sim p_{\text{b}}(y) \end{cases}$

Data dist: $(z, y) \sim p_{\text{data}}$
pairs of images & prompts $\in \mathbb{R}^d$

How would we change our vector field?

Flow models: $u_t^\theta(x) \rightarrow$ Make it prompt dependent

In our case \uparrow guided by prompts.

guided vector field: $u_t^\theta(x|y) \in \mathbb{R}^d$

How to train this? As before

$Z \sim p(\cdot | y)$ $y \rightarrow$ "conditioning" \rightarrow Vector field
 $z \rightarrow$ image of that field
 Conditional \uparrow V.F.
 Prompts in input \rightarrow target

* Guided FM

$L_{CFM}(\theta) = \mathbb{E}_{t, (z, x) \sim p_{data} \times X}$

$\| u_t^\theta(x|y) - u_t(x|z) \|^2$

$t=0 \rightarrow$ Noise
 $t=1 \rightarrow p_{data}(\cdot | y)$

Pairs we are sampling
 (Compare to eqn 26 in notes)

(Show Vanilla guided sampling)

Same way as before except we are sampling a prompt & plugging it into the ODE at every step.

But we get suboptimal results

* How do we make images really fit to the prompt?

Artificially reinforce y

1. Model might underfit
2. text-image pairs may have noise

* Technique: Classifier Free Guidance

But first we talk about classifier guidance

* Classifier Guidance

Bayes Rule

$$P_t(x|y) = \frac{P_t(x) P_t(y|x)}{P_t(y)}$$

Posterior \rightarrow Prior \rightarrow Likelihood \rightarrow Marginal

$$\Rightarrow \log P_t(x|y) = \log P_t(y|x) + \log P_t(x) - \log P_t(y)$$

Applying gradients of log likelihood (the score)

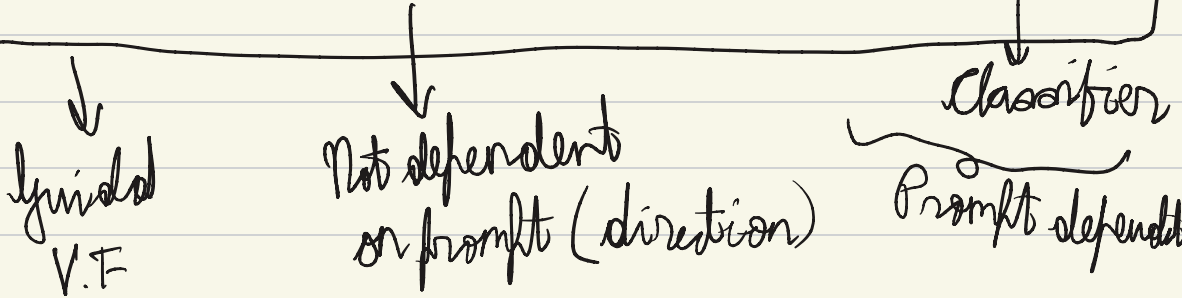
$$\nabla_x \log P_t(x|y) = \nabla_x \log P_t(y|x) + \nabla_x \log P_t(x)$$

The guided score is given by R.H.S. (1)

By using the conversion formula (Eqn 41.8) & substituting ①

$$U_t^{\text{target}}(x|y) = b_t x + a_t \nabla_x \log P_t(x) + \nabla_x \log P_t(y|x)$$

$$U_t^{\text{target}}(x|y) = U_t^{\text{target}}(x) + a_t \nabla_x \log P_t(y|x)$$



$Y \rightarrow$ Label

$x \rightarrow$ Data (Image)

So we scale up the part that is prompt dependent. \rightarrow Show Slides

* Idea: Reinforce Classifier

Weights: $W \geq 1$

$$\tilde{u}_t^w(x|y) = u_t^{\text{target}}(x) + W a_t \nabla \log p_t(y|x)$$

Train the classifier
for guidance.

↓ (add that up)
Scale up contribution

Issues :- 1) We have to train a classifier
along with the flow/diffusion
model.

2) If y is high dimensional,
 $p_t(y|x)$ may be very hard to learn &
 $\nabla \log p_t(y|x)$ is hard to obtain

So Classifier free guidance

* Classifier - free guidance: (CFG)

$$\begin{aligned}
 u_t^w(x|y) &= u_t^{\text{target}}(x) + w a_t \nabla_x \log P_t(y|x) \\
 &= u_t^{\text{target}}(x) + w a_t (\nabla_x \log P_t(x|y) - \nabla_x \log P_t(x)) \\
 &= u_t^{\text{target}}(x) - (w b_t x + w a_t \nabla_x \log P_t(x)) \\
 &\quad + (w b_t x + w a_t \nabla_x \log P_t(x|y)) \\
 &= u_t^{\text{target}}(x) + w (u_t^{\text{target}}(x|y) - u_t^{\text{target}}(x))
 \end{aligned}$$

All of these are vector fields.

$$= w \underbrace{u_t^{\text{target}}(x|y)}_{\text{guided W.F.}} + (1-w) \underbrace{u_t^{\text{target}}(x)}_{\text{unguided W.F.}}$$

$$w > 1$$

There are still
2 models

↓
W.F.

Empty token ϕ : No prompts

Missing prompts

$$\hat{u}_t^w(x|y) = w u_t^\theta(x|y) + (1-w) u_t^\theta(x|\phi)$$

go to the slides up through
this sampling procedure.

* CFG is a heuristic!

We need 2 N.N. calls

1 for conditional

1 for unconditional