

stable Diff from scratch

OUTLINE:

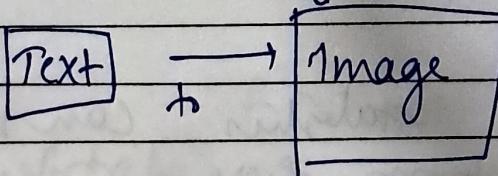
- * Latent Diff Models (stable Diffusion)
- Maths behind it
- Classifier free Guidance
- Text to image.
- Img to img.
- Inpainting -

Prereq:

- (i) Basics of Prob. & statistics (multivariate, Gaussian, conditional prob., marg. prob., Bayes' rule).
- (ii) Basics of Python & NN.
- (iii) Attention Mechanism.
- (iv) How convolutional layers work.

⇒ Stable Diffusion:

↳ introduced by CompVis group in 2022.



Deep Learning Model.

- we will also do img to img (changes in img if req.)
do in images.
- ↳ also inpainting (changes to images).

⇒ Generative Models:

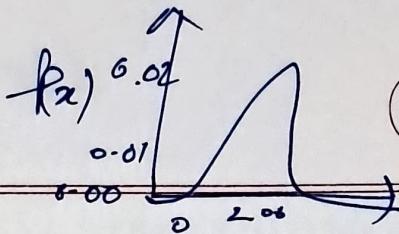
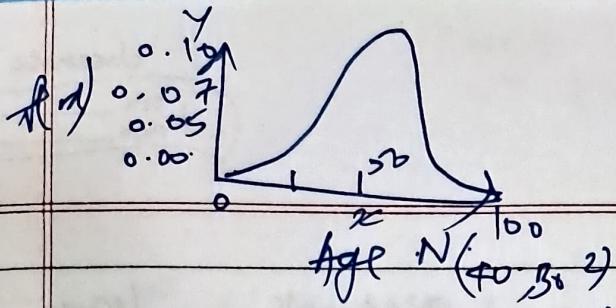
- learn from a probability distribution of the data set such that new instances of data are generated by sampling from it.
- eg: a large sample of cats present from them new samples can be created
- Why do we model data as distributions (probability?)

↳ imagine you are a criminal, trying to create fake identities.

every fake identity is made of variables representing characteristics of person (Age, height).

Age ↳ Height statistics can be used
to sample from these distributions to
create a fake identity.

→ Sample from distrib. means to know a value
from that gives a probability.



classmate

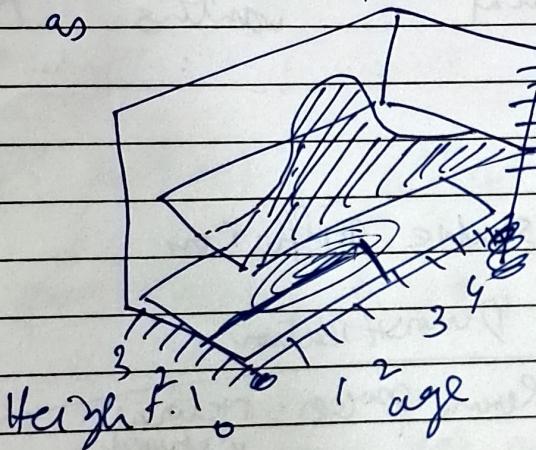
Date _____

Page _____

while producing fake identities let's say.
for eg: assume age = 3, & Height = 130 cm
which is unlikely.

To get reasonable o/p results, we can use joint prob. distributions represented

as



$P(x_1 y)$ represents
the prob. of a certain
 (x, y) combi to occur.

★ To model something we need joint distribution of variables.

★ This is done in image, as well for pixel.

★ Marginalisation: for just Prob. of occur. of x ; regardless of y .

$$P(x) = \int P(x_1 y) P(y) dy$$

$$P(x) = \int p(x, y) dy$$

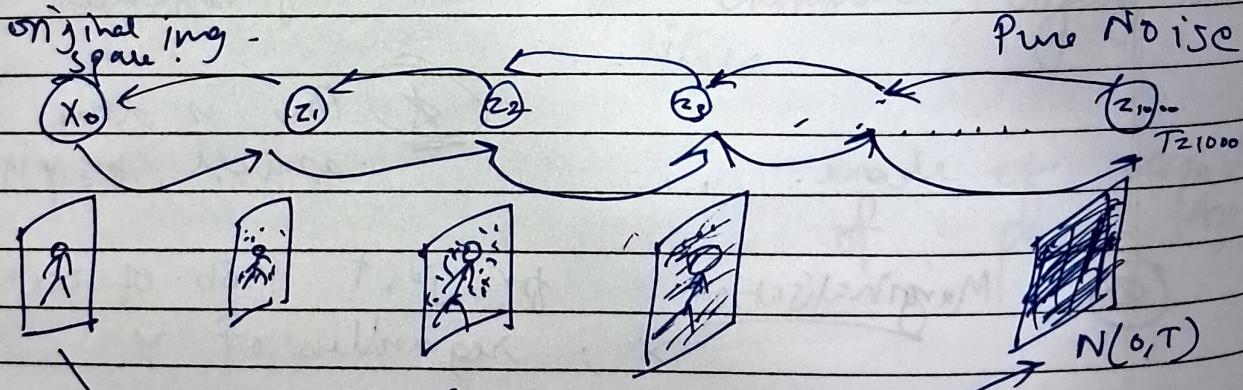
So basically this is what happens in generative models when we learn from prob. distrib. of our data.

- we let the neural network learn the parameters of our distribution.
- we have dataset of images & we want to learn very complex distri that we can't sample from.
just like criminal wanting to generate fake identities.

=) Latent variables : Stable Diffusion

Denoising

Reverse process: neural network



Noising

forward process: fixed

it is easy to add noise
it. than to remove

latent var. models.

classmate

Date _____

Page _____

- Math: we have orig. img how to generate modified image -

$$q(x_t | x_{t-1}) = N(x_t; \underbrace{\sqrt{1-\beta_t}x_{t-1}}_{\text{Mean}}, \underbrace{\beta_t I}_{\text{variance}})$$

β_t = noise int.
Markov chain of var

forward process.

$$q_t(x_t | x_0) = N(x_t; \sqrt{\alpha_t}x_0, (1-\alpha_t)I)$$

from orig. img to img at x_t . without intermediate

noisy \rightarrow less noisy \rightarrow remove noise.
new network learns parameters. how to change

$$p_\theta(x_{t-1} | x_t) = N(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

\Rightarrow Diffusion Models: $p_\theta(x_0) = \int p_\theta(x_0; \tau) dx_{\cdot \tau}$.

\Rightarrow How do we train?

Algorithm:

1. repeat.
2. $x_0 \sim q(x_0)$
3. $t \sim \text{Uniform}(\{1, \dots, T\})$
4. $\epsilon \sim \mathcal{N}(0, I)$
5. g_t depend on D_0 || $e^{-\epsilon_t} (1 - \sqrt{1 - \epsilon_t})$

6. until converged / min.

⇒ what we do is that we optimise the lower bound like for eg:

$$\text{Revenue} \geq \text{Sales}$$

lower bound

to inc. Revenue ↑ inc. Sales ↑.

Same for us.

⇒ Based on Algo $E_0 \rightarrow NN$.

$$E_0(\sqrt{\bar{\alpha}_t x_t} + \sqrt{1-\bar{\alpha}_t} \epsilon_t)$$

Noisy img time t at which noise added

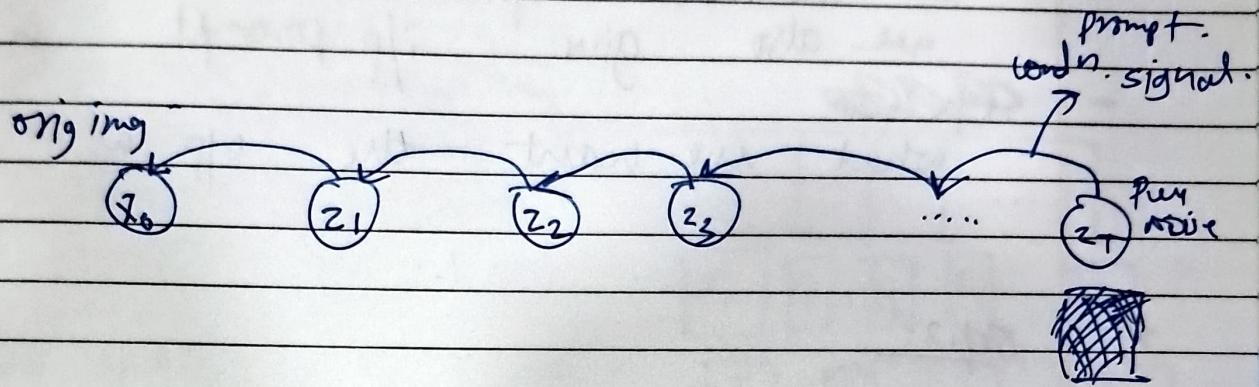
NN has to calc. how much noise added / prevent. On applying GD, we maximise ELBO, & log likelihood.

⇒ Sampling New data from noisy images at T time

- During the Rev. process we recursively calc. how much noise there is & remove it by help of NN.

- by remov. noise from pure noise we generate new data.

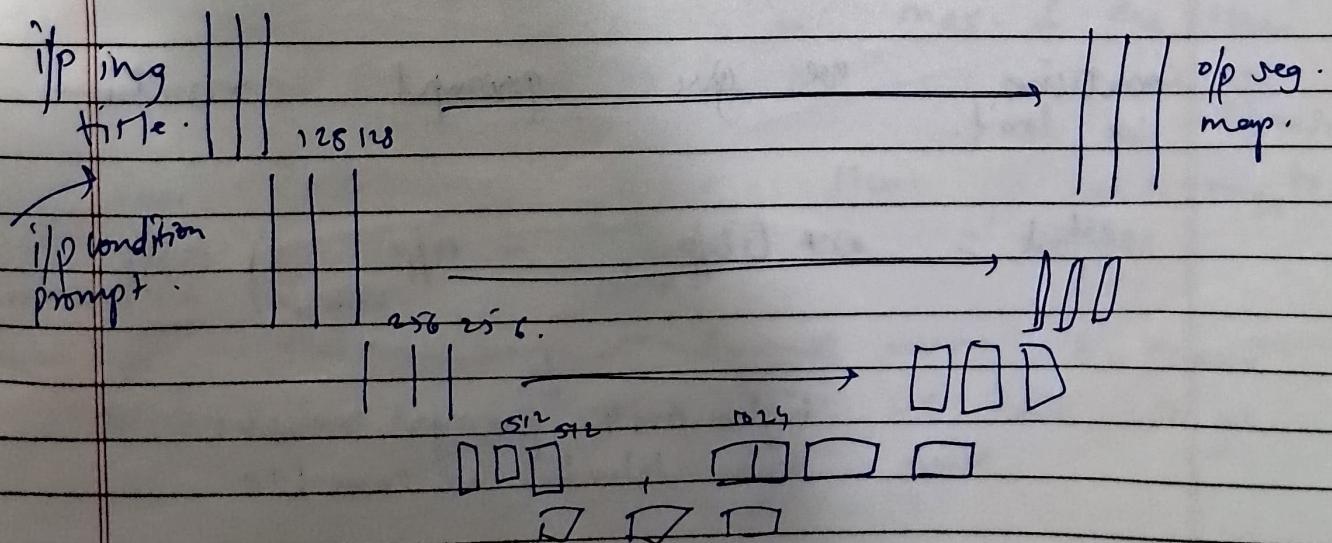
- But to create something we want & sensible we want to control the misification process.



- Our final goal is to model a distrib. Theta $P(\cdot)$ such that $P_\theta(x)$ is such that we maximize likelihood of our data & learn this distrib. we maximize the elbo (lower bound) by minimize the loss ∇_θ which learns the distrib.

- we use the vNet Model which receives a noisified img and has to find amt of noise.

16464



* - Step 1:

give UNET model noised image $T=1000$

& the model predicts the noise level.

we also give i/p prompt as to

~~what we want~~ what we want the o/p as CAT etc

* - Step 2:

we give same i/p as noise but now
we don't give i/p. prompt, so the
model will build some o/p to generate
something.

- This way we combine o/p of the both
steps to get type the ~~to~~ type
of output we want closer to prompt

classifier free guidance.

- Sometimes we give prompt sometimes
we don't.

$$\text{output} = w * (\text{o/p condn.} - \underbrace{\text{o/p uncondn.}}_{\substack{\text{higher this} \\ \downarrow \\ \text{prompt resembles}}} + \text{o/p uncondn.})$$

higher this prompt resembles.
lower value less resemble.

→ Slow based models.
↳ classifier Grid.

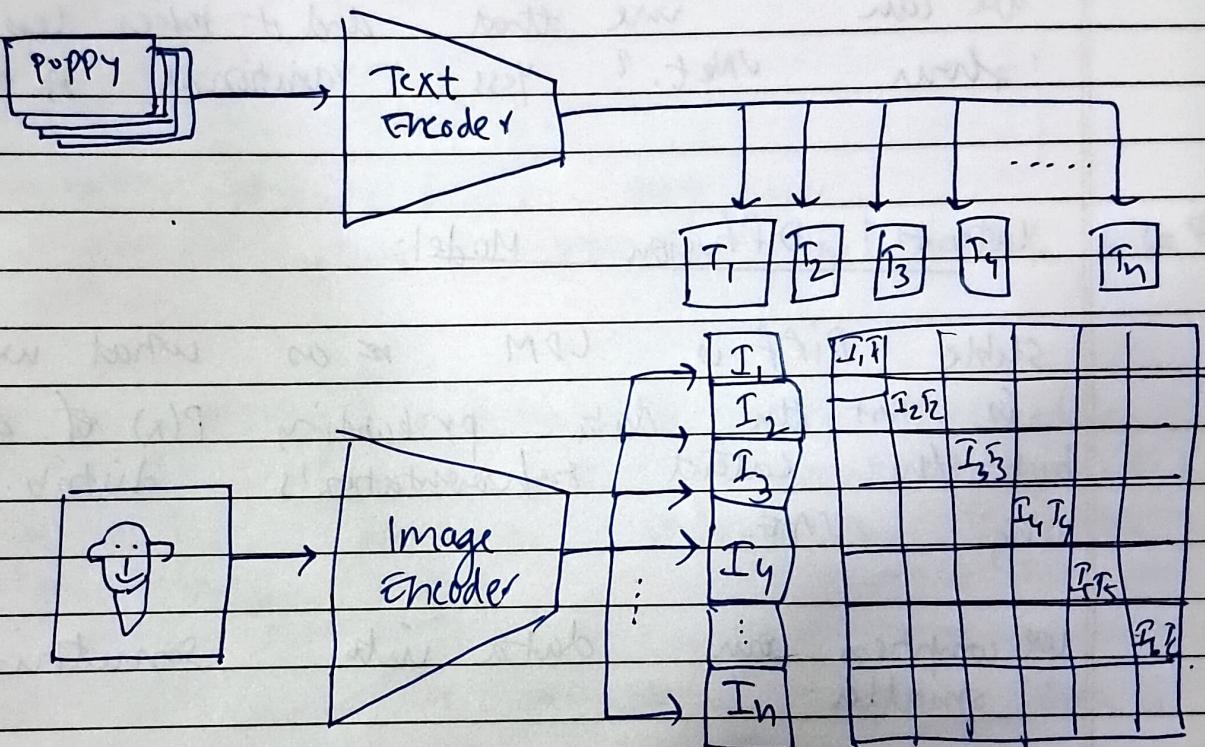
classmate

Date _____

Page _____

⇒ CLIP (contrastive language - image pre-training):

1.) Contrastive Pretraining: (Open AI)



- we create a loss function where the diagonal must I_{ii} must be max. & the other values of matrices to be 0.

signals.

These embeddings act as the i/p conditions to the model. to denoise the image.

- Now the thing is it is very big process to do repeatedly every time we have to go to UNET for finding noise.

- If img is big (512×512) then big matrix will be thru.

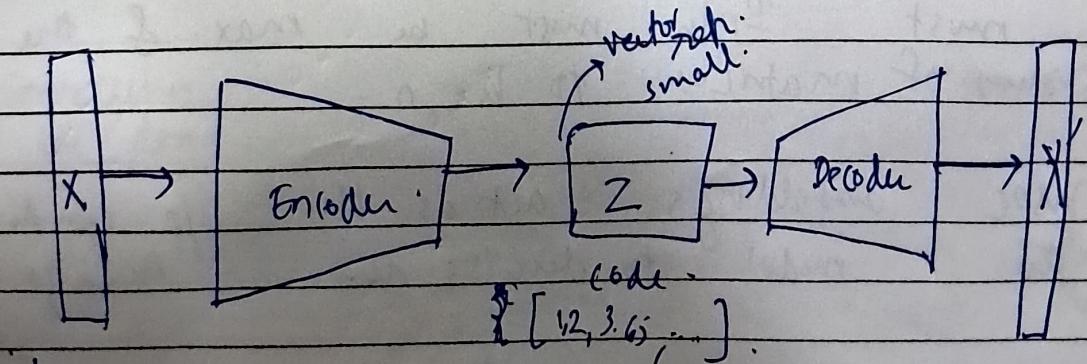
can we compress the image so that we can use that and it takes less time than vNet.? Yes variational Auto Encoder.

\Rightarrow Latent Diffusion Model:

stable Diff is LDM \Rightarrow as what we learn is not the data probability $P(x)$ of dataset but the latent representation's distrib. by using VAE.

We compress our data into something smaller.

It is like zipping a large file : Something same.



i/p.

Reconstructed
i/p.

- Orig image passes thru the Encoder and is conv. to vector format resized & then passing thru the decoder convert as small. but same decodin. i/p.

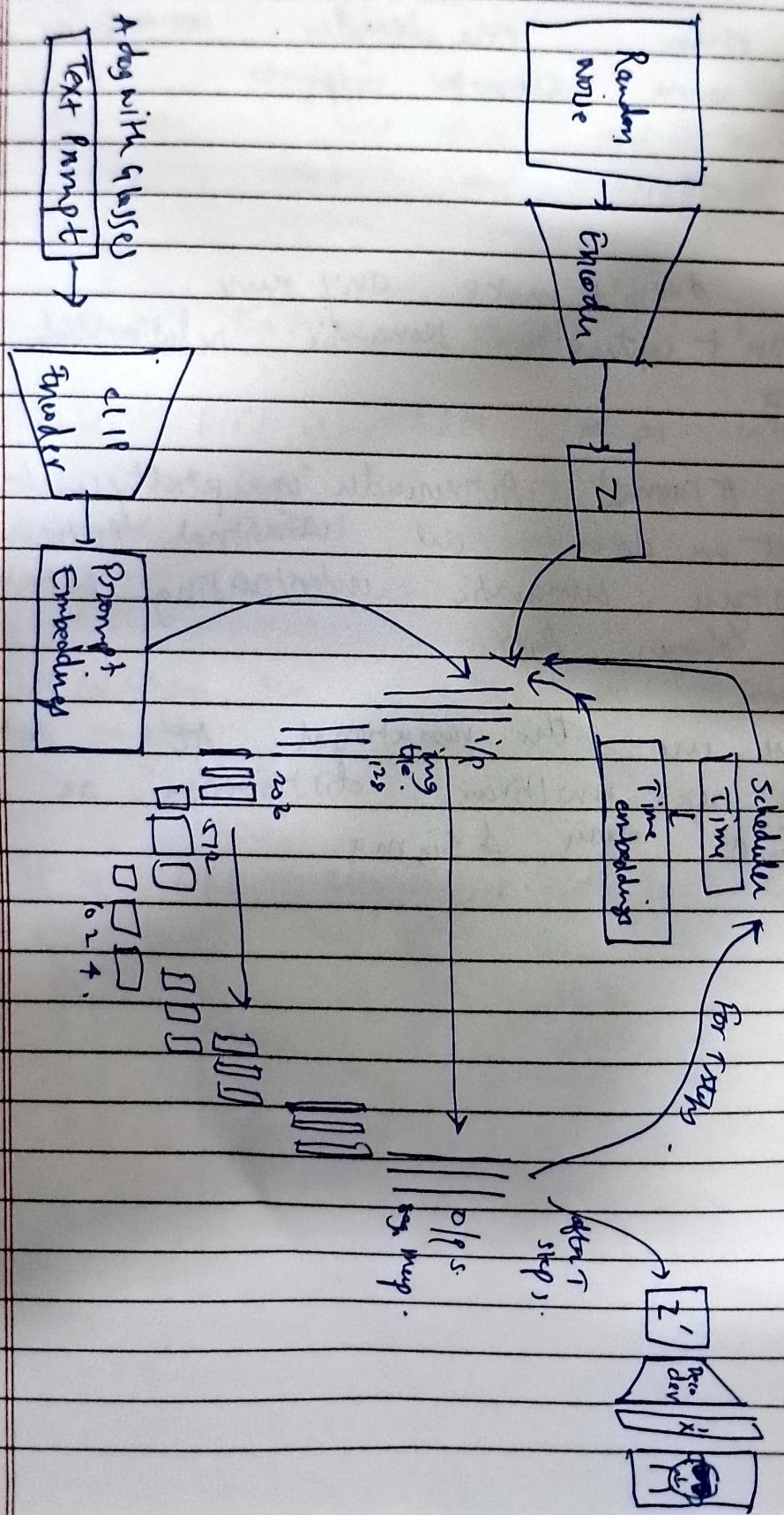
- Issues:

- wde doesn't make any sense.
- doesn't capture semantic relationship b/w data.

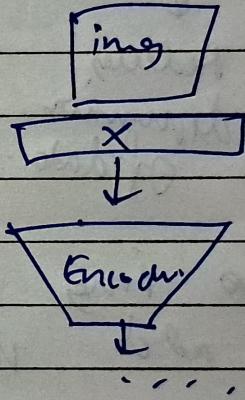
NOTE: A normal Autoencoder is pretty useless, what we do is call variational Autoencoder to capture semantic understanding between the ~~changes~~ data.

Hence we use the variational AE. data is diffi across multivar. distribution, as a gaussian & find mean & sigma.

= Architecture of stable Diffusion: (Text to image)



- in case of Img-to-image, we input img.
- * and then provide IIP for any class to be made. ∴



& the Rest is same.

We noisy the img & the point that we noisy demonstrates now much freedom user has to modify img.