

INTRODUCTION TO
MACHINE LEARNING
COMPSCI 4ML3

LECTURE 28

HASSAN ASHTIANI

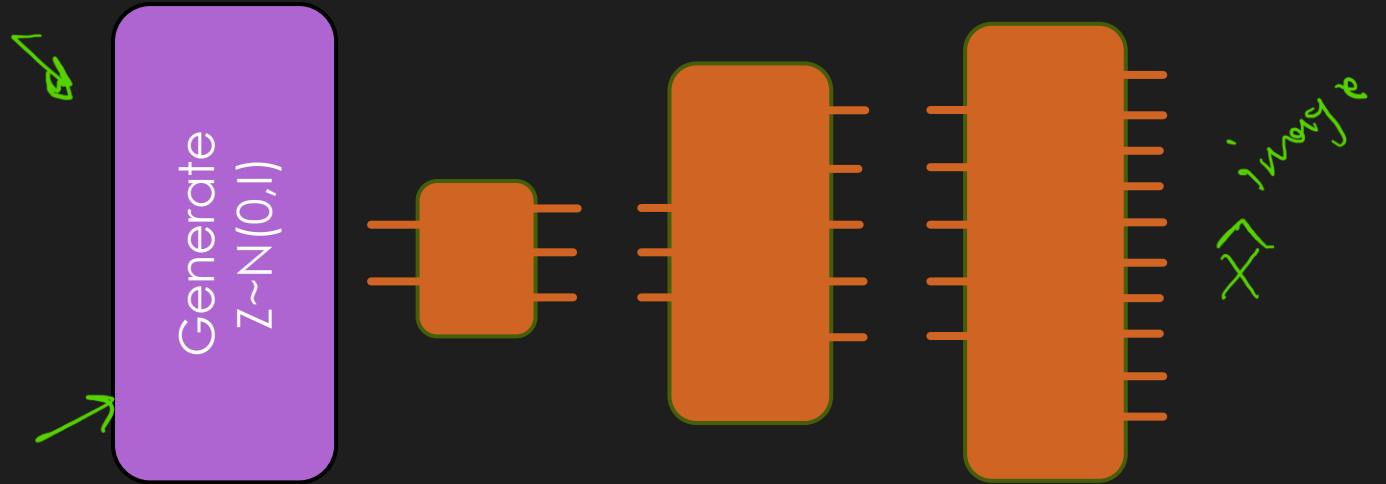
MAXIMIZE LIKELIHOOD?

$$P[\hat{x}] = \text{[sad face]}$$

- $\text{MAX}_{\theta} \sum_x \widehat{p(x)}$
(log L)

$\rightarrow p(x|z)$: easy ✓
 $p(x)$ x hard

$\int_z p(x|z) p(z) = p(x)$
integration is slow...



TEXT TO IMAGE MODELS

- SUCCESS DUE TO: LLMs + DIFFUSION MODELS



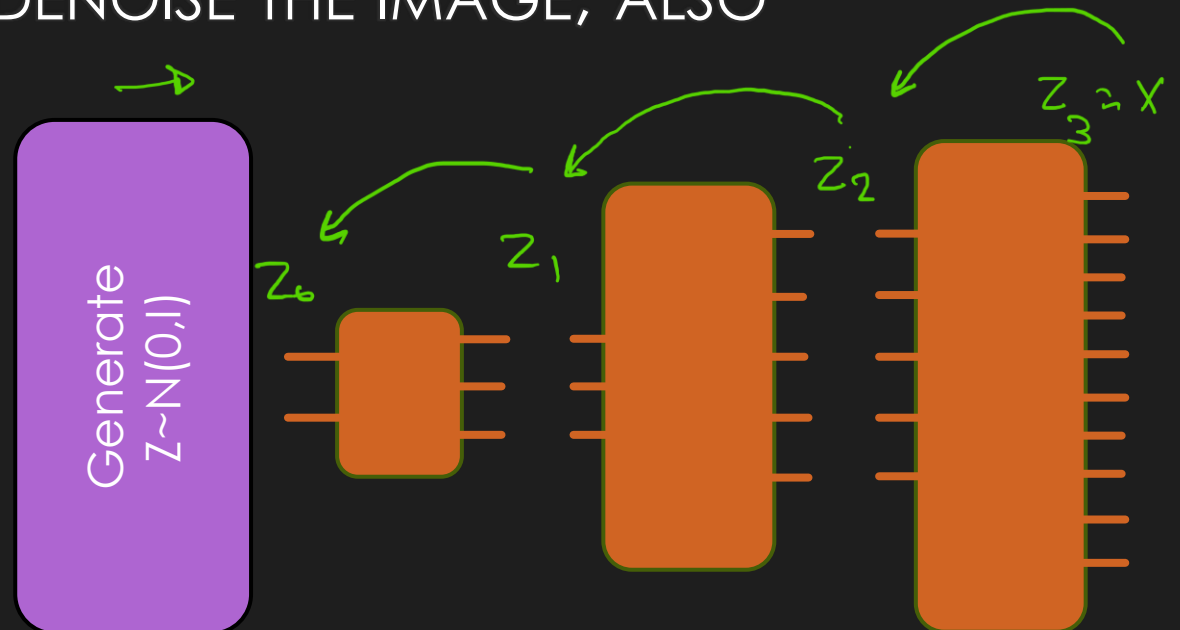
DALL-E 2



IMAGEN

DIFFUSION MODELS

- VIEW THE NETWORK IN TWO DIRECTIONS:
 - ➔ • RIGHT TO LEFT: PROGRESSIVELY ADD MORE NOISE, ALSO MAKE THE IMAGE SMALLER IN RESOLUTION
 - ➔ • LEFT TO RIGHT: PROGRESSIVELY DENOISE THE IMAGE, ALSO MAKE IT HIGH RESOLUTION



TRAINING DIFFUSION MODELS

- $t \in [T]$ REPRESENTS A LAYER OF THE NETWORK
- PROGRESSIVELY ADDING NOISE TO AN IMAGE:

- $z_T = x$
- $z_{t-1} = \text{Noisy}(z_t)$ FOR ALL t
- $z_0 = \text{Noisy}(z_1)$

- PROGRESSIVELY DENOISING AN IMAGE:

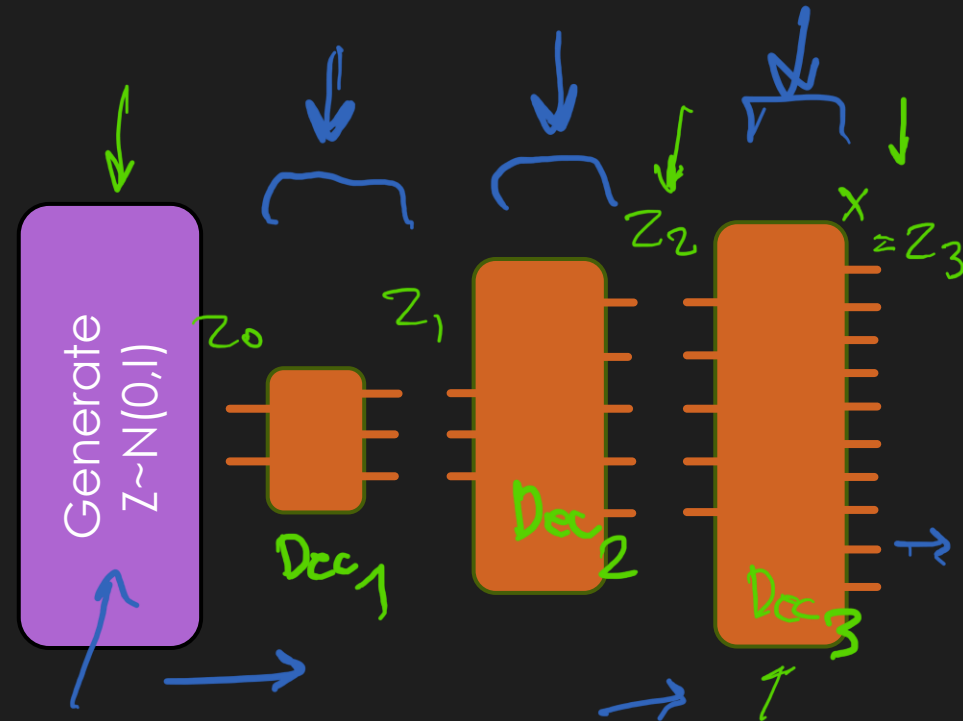
- $\hat{z}_t = \text{Dec}_t(z_{t-1})$

- GENERATE AN IMAGE FROM $z_0 \sim N(0, I)$

- $\hat{x} = \text{Dec}(z_0) = \text{Dec}_T(\dots \text{Dec}_2(\text{Dec}_1(z_0)) \dots)$

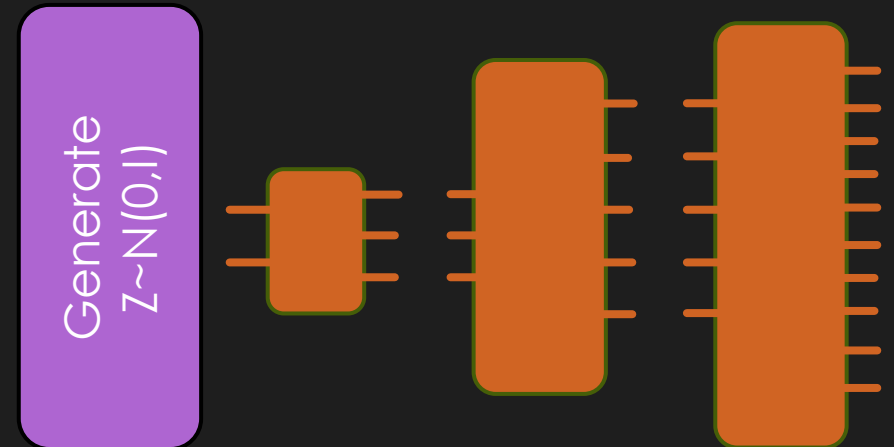
- TRAINING DENOISERS

loss for one data point $\sum_t w_t \| \text{Dec}_t(z_{t-1}) - z_t \|_2^2$



CONDITIONING IMAGE GENERATION

- HOW TO CONTROL WHAT THE MODEL GENERATES?
 - E.G., FORCE TO GENERATED IMAGE FROM CLASS c ?
 - ASSUME WE HAVE LABELED DATA (x, c)
- DENOISER RECEIVES LABEL AS WELL
 - $\hat{z}_t = Dec_t(z_{t-1}, c)$
- TRAINING DENOISERS
 - $\sum_{(x,c)} w_t \|Dec_t(z_{t-1}, c) - z_t\|_2^2$
- IMAGEN: RANDOMLY DROP c 10% OF THE TIMES DURING TRAINING
- FOR TEXT TO IMAGE MODELS, c CAN BE THE LATENT REPRESENTATION OF TEXT RATHER THAN JUST THE LABEL.



ADVERSARIAL PERTURBATIONS

TYPICAL CLASSIFIERS ARE SENSITIVE TO (IMPERCEPTIBLE)
“ADVERSARIAL” PERTURBATIONS

SZEGEDY ET AL.'14



RISKS OF ADVERSARIAL PERTURBATIONS

- VULNERABLE TO MALICIOUS ATTACKS
- IGNORE THE “INVARIANCE”/DOMAIN-KNOWLEDGE
- WHY ARE THESE CALLED ADVERSARIAL?
 - THE NOISE IS NOT RANDOM
 - CAREFULLY SELECTED TO FOOL THE MODEL



DECISION BOUNDARY VISUALIZATION

- SMOOTHING THE DECISION BOUNDARY HELPS



FINDING ADVERSARIAL PERTURBATIONS

- USUAL GRADIENT DESCENT TO FIND PARAMS θ

- $\nabla_{\theta}(\underline{L}(\theta)) = \frac{1}{m} \sum_i \nabla_{\theta} (l(f_{\theta}(x^i), y^i))$

- $\underline{\theta} = \underline{\theta} - \alpha \nabla_{\theta}(\underline{L}(\theta))$

- GRADIENT ASCENT FOR FINDING PERTURBATIONS

- $\nabla_x(l(x, y)) = \nabla_x (l(f_{\theta}(x^i), y^i))$

- $x = x + \alpha \nabla_x(l(x, y))$

