

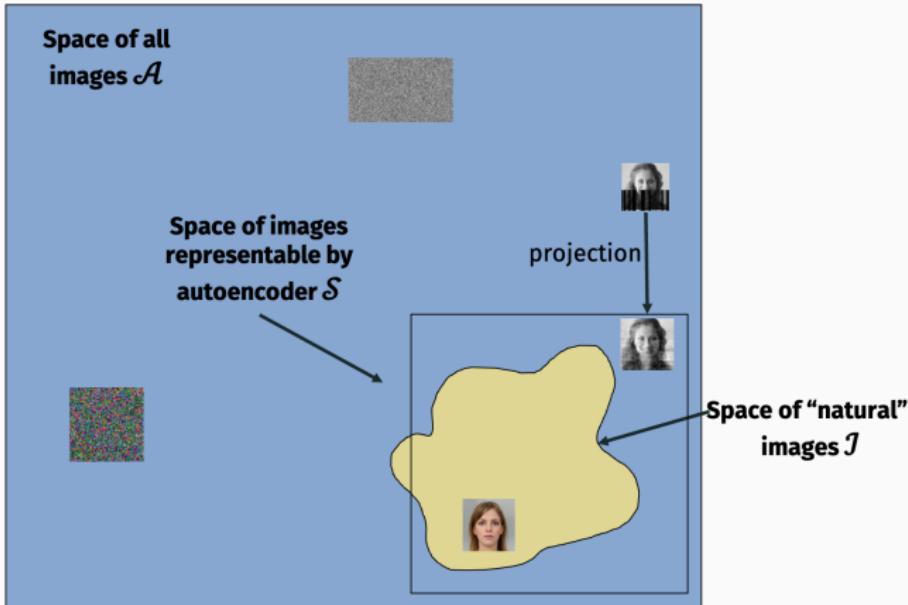
# **CS-GY 6923: Lecture 14**

## **Image Generation, and Privacy Concerns in ML**

---

NYU Tandon School of Engineering, Akbar Rafiey

## Recap: Autoencoders learn compressed representations



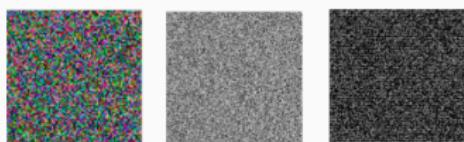
$f(\mathbf{x}) = d(e(\mathbf{x}))$  projects an image  $\mathbf{x}$  closer to the space of natural images.

## Autoencoders for data generation

Suppose we want to generate a random natural image. How might we do that?

- **Option 1:** Draw each pixel value in  $\mathbf{x}$  uniformly at random.

Draws a random image from  $\mathcal{A}$ .



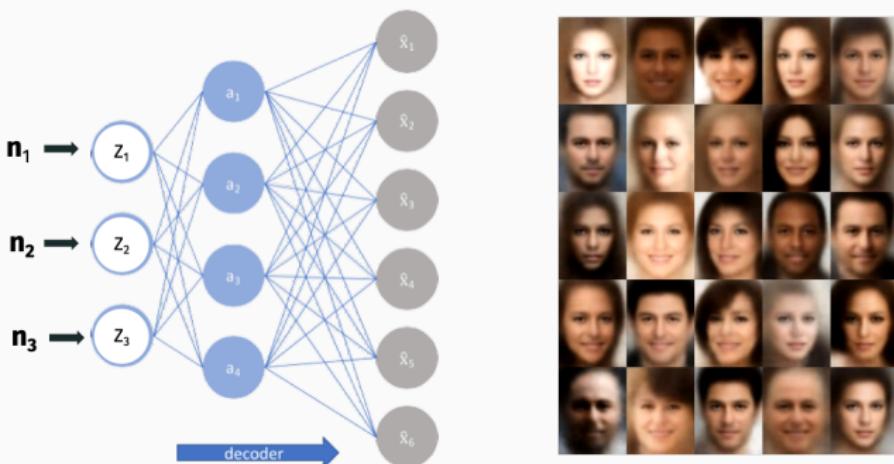
- **Option 2:** Draw  $\mathbf{x}$  randomly from  $\mathcal{S}$ , the space of images representable by the autoencoder.



How do we randomly select an image from  $\mathcal{S}$ ?

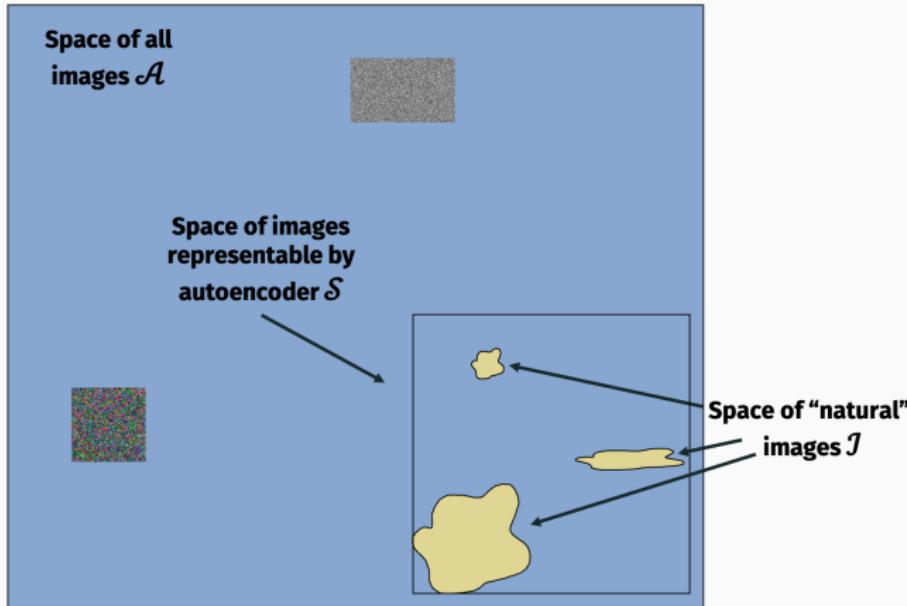
# Autoencoders for data generation

**Autoencoder approach to generative ML:** Feed random inputs into decoder to produce random realistic outputs.



**Main issue:** most random inputs will “miss” and produce garbage results.

# Autoencoders for data generation



Variational Auto-Encoders (VAEs) attempt to resolve this issue.

# Variational AutoEncoders (VAEs)

VAEs attempt to resolve this issue. Basic ideas:

- Instead of mapping inputs to a single latent vector, VAEs map them to a probability distribution in the latent space (e.g., a Gaussian distribution)

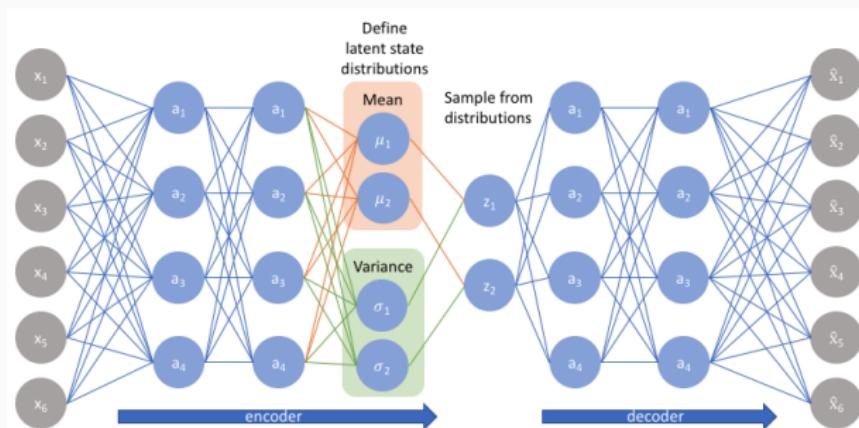


Image from <https://www.jeremyjordan.me/variational-autoencoders/>

# Variational AutoEncoders (VAEs)

Basic ideas:

- Suppose there exists some hidden variable  $z$  which generates  $x$ .
- Ideally we want to understand  $p(z | x)$  (probabilistic encoder) and  $p(x | z)$  (probabilistic decoder)
- We can only see the data. Computing  $p(z | x)$  is hard

$$p(z | x) = \frac{p(x | z)p(z)}{p(x)}$$

# Variational AutoEncoders (VAEs)

Basic ideas:

- Let's approximate  $p(z | x)$  using a simpler to understand distribution  $q(z | x)$ .

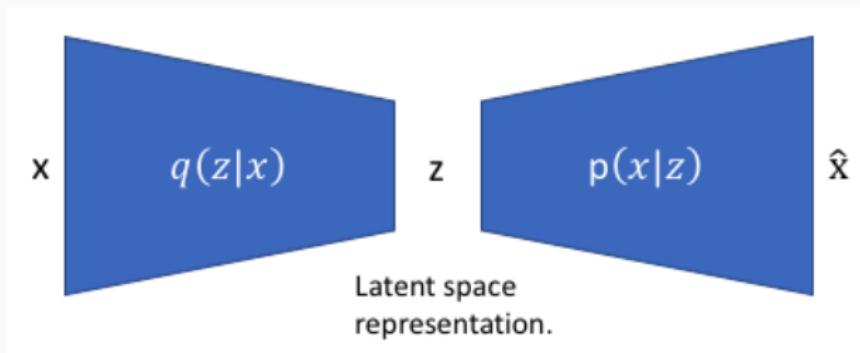


Image from <https://www.jeremyjordan.me/variational-autoencoders/>

# Variational AutoEncoders (VAEs)

Basic ideas:

- Let's approximate  $p(z | x)$  using a simpler to understand distribution  $q(z | x)$ .
- $q(z | x)$  must have nice properties e.g., it should be as similar as possible to  $p(z | x)$
- Optimization problem !

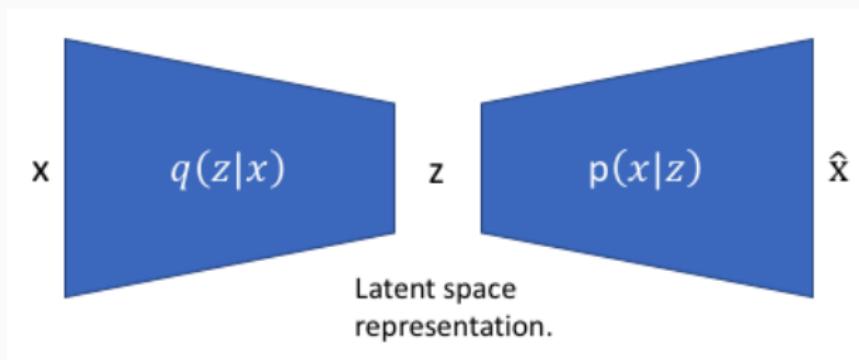


Image from <https://www.jeremyjordan.me/variational-autoencoders/>

## Kullback–Leibler divergence

KL divergence is a measure of difference between two probability distributions.

$$KL(P, Q) = \sum_x P(x) \frac{\log(P(x))}{\log Q(x)}$$

Back to our optimization problem:  $q(z | x)$  and  $p(z | x)$  should be similar

$$\min KL(q(z | x), p(z | x))$$

Equivalent to:

$$\max \mathbb{E}_{q(z|x)} \log(p(x | z)) - KL(q(z | x), p(z))$$

## VAE objective

Back to our optimization problem:  $q(z | x)$  and  $p(z | x)$  should be similar

$$\min KL(q(z | x), p(z | x))$$

Equivalent to:

$$\max \mathbb{E}_{q(z|x)} \log(p(x | z)) - KL(q(z | x), p(z))$$

What is the first term? what is the second term?

## VAE objective

Back to our optimization problem:  $q(z | x)$  and  $p(z | x)$  should be similar

$$\min KL(q(z | x), p(z | x))$$

Equivalent to:

$$\max \mathbb{E}_{q(z|x)} \log(p(x | z)) - KL(q(z | x), p(z))$$

First term: reconstruction likelihood.

Second term: ensures that our learned distribution  $q$  is similar to the true prior distribution  $p$ .

# VAEs: implementation

Have neural networks to learn the mappings  $q(z | x)$  and  $p(x | z)$ .

$$\max_{\theta, \phi} \mathbb{E}_{q_\phi(z|x)} \log(p_\theta(x|z)) - KL(q_\phi(z|x), p_\theta(z))$$

Assumptions:  $q_\phi(z|x)$ ,  $p_\theta(x|z)$ , and  $p_\theta(z)$  are Gaussian distributions.

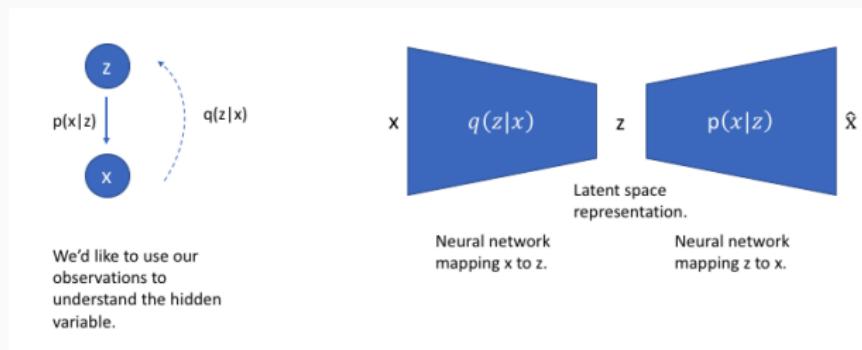


Image from <https://www.jeremyjordan.me/variational-autoencoders/>

# VAEs: implementation

The encoder model of a VAE will output parameters describing a distribution for each dimension in the latent space.

Assuming Gaussian we only need a mean and a variance for describing each dimension in the latent space

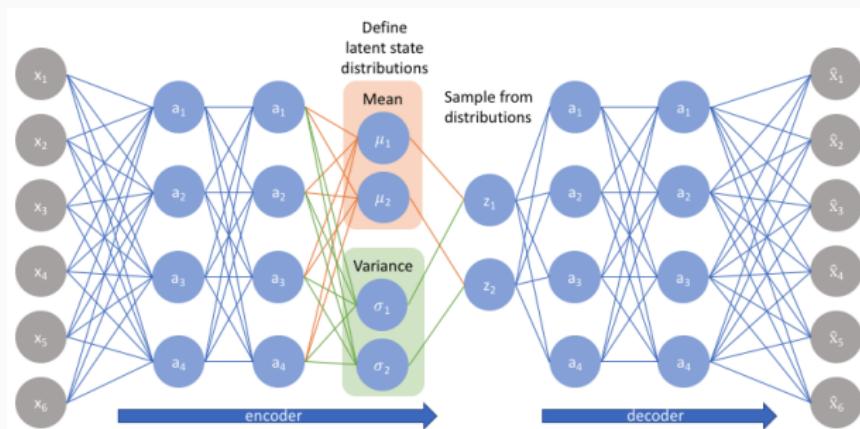


Image from <https://www.jeremyjordan.me/variational-autoencoders/>

# VAEs: implementation

It is not easy to backpropagate the gradient through samples !

**Reparameterization** technique addresses this issue.

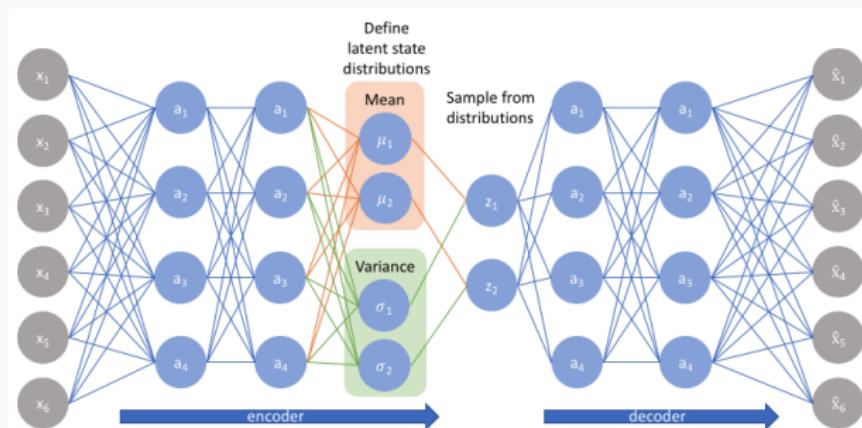


Image from <https://www.jeremyjordan.me/variational-autoencoders/>

## VAEs: reparameterization technique

$\mathbf{z} \sim q_\phi(\cdot | \mathbf{x})$  is normally distributed, as

$$\mathcal{N}(\mu_\phi(\mathbf{x}), \sigma_\phi(\mathbf{x}))$$

This can be reparameterized by letting  $\varepsilon \sim \mathcal{N}(0, I)$ , and constructing  $\mathbf{z}$  as

$$\mathbf{z} = \mu_\phi(\mathbf{x}) + L_\phi(\mathbf{x})\varepsilon$$

Here,  $\sigma_\phi(\mathbf{x})$  is obtained by the Cholesky decomposition:

$$\sigma_\phi(\mathbf{x}) = L_\phi(\mathbf{x})L_\phi(\mathbf{x})^T$$

Then we have

$$\nabla_\phi \mathbb{E}_{\mathbf{z} \sim q_\phi(\cdot | \mathbf{x})} \left[ \ln \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x})} \right] = \mathbb{E}_\varepsilon \left[ \nabla_\phi \ln \frac{p_\theta(\mathbf{x}, \mu_\phi(\mathbf{x}) + L_\phi(\mathbf{x})\varepsilon)}{q_\phi(\mu_\phi(\mathbf{x}) + L_\phi(\mathbf{x})\varepsilon | \mathbf{x})} \right]$$

and so we obtain an unbiased estimator of the gradient, allowing stochastic gradient descent.

## VAEs: reparameterization technique

Since we reparameterized  $z$ , we need to find  $q_\phi(z|x)$ . Let  $q_0$  be the probability density function for  $\varepsilon$ , then

$$\ln q_\phi(z|x) = \ln q_0(\varepsilon) - \ln |\det(J(z, \varepsilon))|$$

Since  $z = \mu_\phi(x) + L_\phi(x)\varepsilon$ , this becomes

$$\ln q_\phi(z|x) = -\frac{1}{2}\|\varepsilon\|^2 - \ln |\det L_\phi(x)| - \frac{n}{2} \ln(2\pi)$$

## VAEs: reparameterization technique

We can now optimize the parameters of the distribution (unbiased estimation of the gradient) while still maintaining the ability to randomly sample from that distribution.

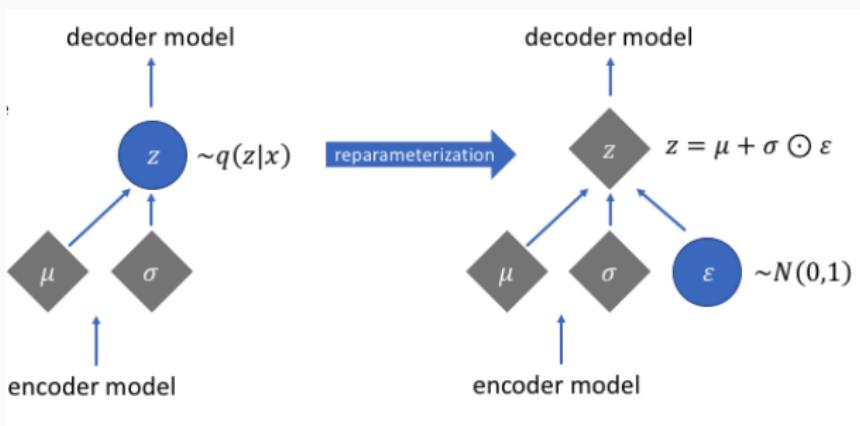
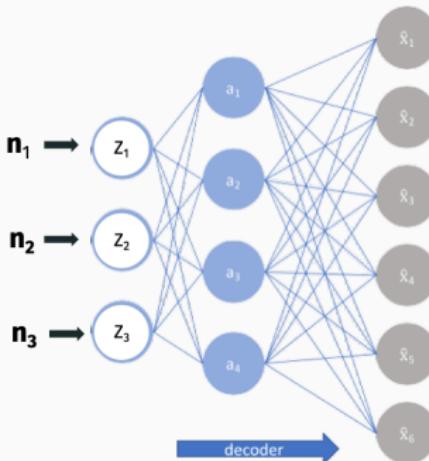


Image from <https://www.jeremyjordan.me/variational-autoencoders/>

# Generative Adversarial Networks (GANs)

VAEs give very good results, but tends to produce images with immediately recognizable flaws (e.g. soft edges, high-frequency artifacts).



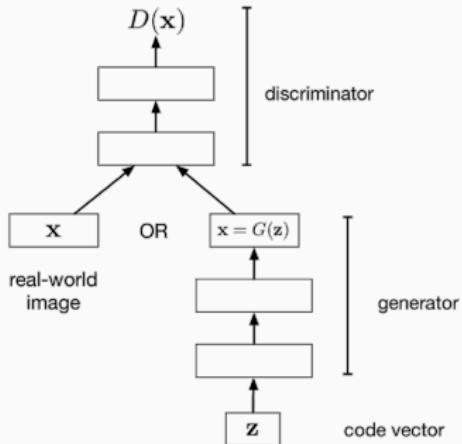
# Generative Adversarial Networks (GANs)

Lots of efforts to hand-design regularizers that penalize images that don't look realistic to the human eye.

**Main idea behind GANs:** Use machine learning to automatically encourage realistic looking images.

$$\min_{\theta} L(\theta) + P(\theta)$$

# Generative Adversarial Networks (GANs)

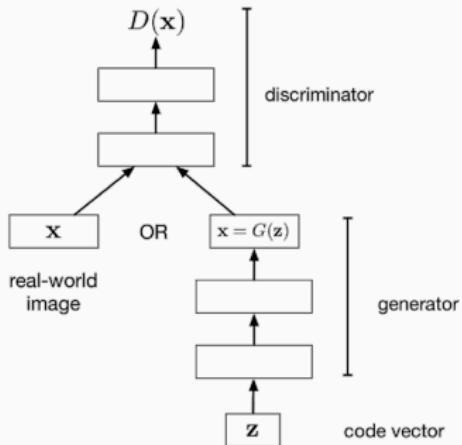


Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be real images and let  $\mathbf{z}_1, \dots, \mathbf{z}_m$  be random code vectors. The goal of the discriminator is to output a number between  $[0, 1]$  which is close to 0 if the image is fake, close to 1 if it's real.

Train weights of discriminator  $D_\theta$  to minimize:

$$\min_{\theta} \sum_{i=1}^n -\log(D_\theta(\mathbf{x}_i)) + \sum_{i=1}^m -\log(1 - D_\theta(G_\theta(\mathbf{z}_i)))$$

# Generative Adversarial Networks (GANs)



Goal of the generator  $G_{\theta'}$  is the opposite. We want to maximize:

$$\max_{\theta'} \sum_{i=1}^m -\log (1 - D_{\theta}(G_{\theta'}(z_i)))$$

This is called an “adversarial loss function”.  $D$  is playing the role of the adversary.

# Generative Adversarial Networks (GANs)

$$\theta^*, \theta'^* \text{ solve } \min_{\theta} \max_{\theta'} \sum_{i=1}^n -\log(D_{\theta}(\mathbf{x}_i)) + \sum_{i=1}^m -\log(1 - D_{\theta}(G_{\theta'}(\mathbf{z}_i)))$$

This is called a minimax optimization problem. Really tricky to solve in practice.

- **Repeatedly play:** Fix one of  $\theta^*$  or  $\theta'^*$ , train the other to convergence, repeat.
- **Simultaneous gradient descent:** Run a single gradient descent step for each of  $\theta^*, \theta'^*$  and update  $D$  and  $G$  accordingly. Difficult to balance learning rates.
- Lots of tricks (e.g. slight different loss functions) can help.

# Generative Adversarial Networks (GANs)

State of the art until a few years ago.



## Quality of generative model

How to evaluate the quality of our model e.g., GAN? What do we expect from it?

## Quality of generative model

How to evaluate the quality of our model e.g., GAN? What do we expect from it?

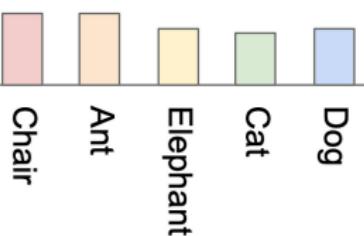
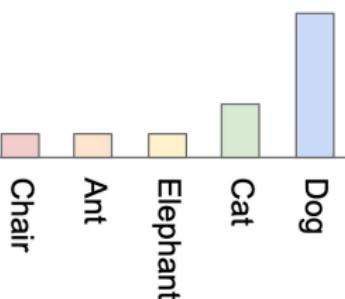
- The images generated by our model should have variety (e.g., each image is a different breed of dog)
- Each image distinctly looks like something (e.g., one image is clearly a Poodle, the next a great example of a French Bulldog)
- ...

# The Inception Score (IS)

Each image distinctly looks like something (e.g., one image is clearly a Poodle, the next a great example of a French Bulldog)



Trained Model:  
gives a  
distribution  
over classes

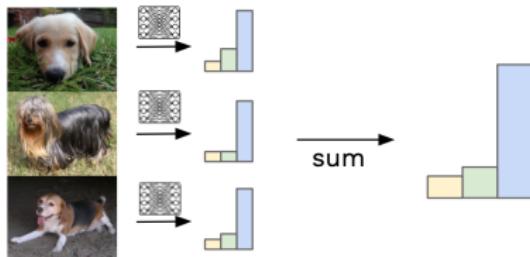


# The Inception Score (IS)

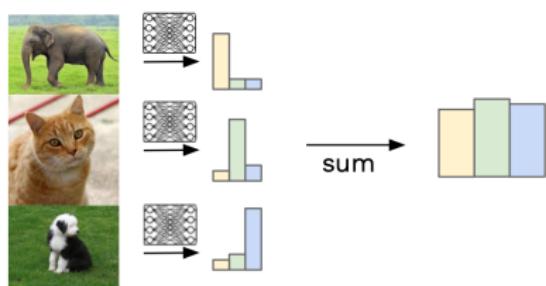
The images generated by our model should have variety:

Generate a lot of images (50,000) using the model and sum their distributions.

Similar labels sum to give focussed distribution



Different labels sum to give uniform distribution



# The Inception Score (IS)

Higher KL divergence, means better score.

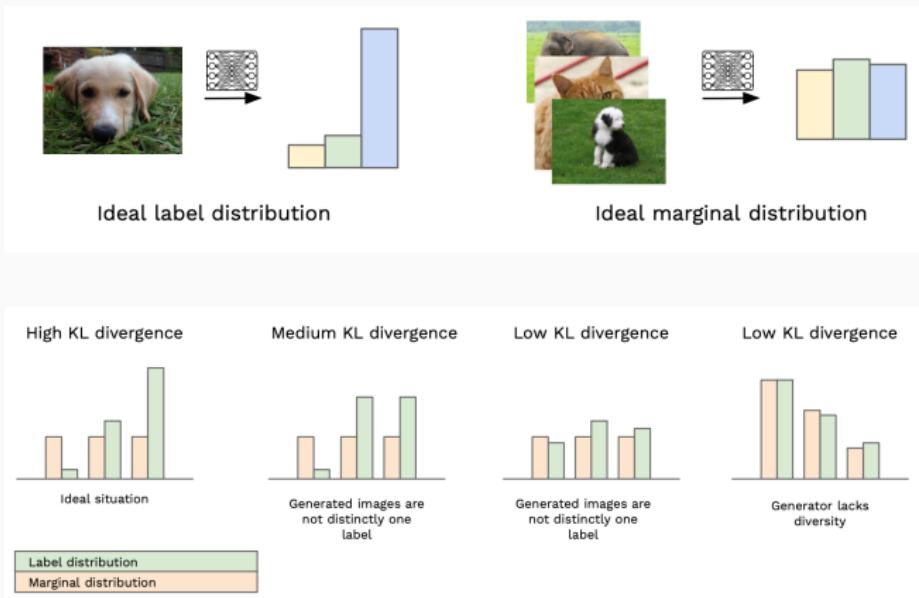


Image from

<https://medium.com/octavian-ai/a-simple-explanation-of-the-inception-score-372dff6a8c7a>

# Diffusion Model

**Auto-encoder/VAE, GAN approach:** Input noise, map directly to image and vice versa.

**Diffusion:** Slowly move from noise to image and vice versa.

---

## Denoising Diffusion Probabilistic Models

---

**Jonathan Ho**

UC Berkeley

[jonathanho@berkeley.edu](mailto:jonathanho@berkeley.edu)

**Ajay Jain**

UC Berkeley

[ajayj@berkeley.edu](mailto:ajayj@berkeley.edu)

**Pieter Abbeel**

UC Berkeley

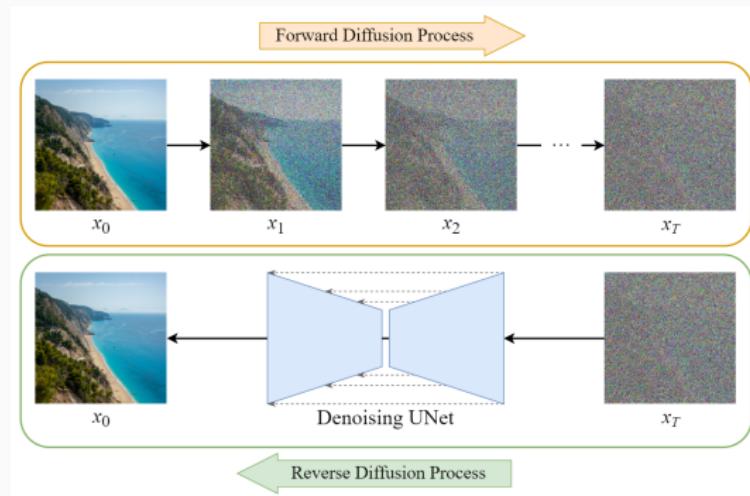
[pabbeel@cs.berkeley.edu](mailto:pabbeel@cs.berkeley.edu)

### Abstract

We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic

# How diffusion models work

- Forward Process:
  - Gradually add noise to data until it becomes pure noise.
- Reverse Process:
  - Train a neural network to remove the noise step by step.

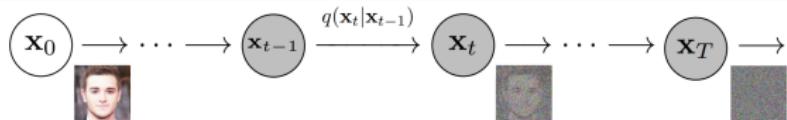


Key Question: How do we predict and reverse noise effectively?

# Mathematical Formulation (1/2)

## Forward Process (Adding Noise):

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$



- $\beta_t$ : Noise schedule.
- After  $T$  steps, for large enough  $T$ ,  $\mathbf{x}_T$  is pure noise.

## Cumulative Noise:

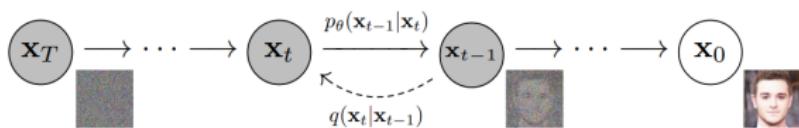
$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

with retention factor  $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$ .

## Mathematical formulation (2/2)

Reverse Process (Denoising):

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$$



- $\mu_\theta$ : Predicted mean of the clean image.
- $\Sigma_\theta$ : Predicted variance (optional).

Training objective:

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0, t, \epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2]$$

## Diffusion models vs VAEs

- Recall the goal of VAE was to have a probabilistic representation of latent space attributes. This was done in one shot, from image to Gaussian distributions.
- VAE decoder does the reverse in one shot. Takes Gaussian noise and returns an image.
- Diffusion model doing more or less the same but with many careful intermediate noising and denoising steps.
- There are fundamental differences ...

## Diffusion process

- A diffusion process is a stochastic Markov process having continuous path
- stochastic Markov process: future state will only depend on the current state, knowing the past does not change anything
- continuous: no jumps
- Allows to transition from complex distributions to simple distributions

## Forward process: a closer look

Forward process:  $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$

Why does it converge to a simple distribution ?

Suppose the process is  $\mathbf{x}_t = \alpha \mathbf{x}_{t-1} + \beta \mathcal{N}(0, I)$ .

What values of  $\alpha$  and  $\beta$  makes sense for the process?

- $\alpha = 0, \beta = 1$ ?
- $\alpha \geq 1, \beta \leq 1$ ?
- Seems something like  $\alpha = \sqrt{0.99}, \beta = \sqrt{0.01}$  is appropriate.

## Forward process: a closer look

Let's check if the following  $\alpha$  and  $\beta$  converges:

$$\mathbf{x}_t = \sqrt{1 - \beta} \mathbf{x}_{t-1} + \sqrt{\beta} \mathcal{N}(0, I)$$

$$\begin{aligned}\mathbf{x}_t &= \sqrt{1 - \beta} \mathbf{x}_{t-1} + \sqrt{\beta} \mathcal{N}(0, I) \\ &= \sqrt{1 - \beta} (\sqrt{1 - \beta} \mathbf{x}_{t-2} + \sqrt{\beta} \mathcal{N}(0, I)) + \sqrt{\beta} \mathcal{N}(0, I) \\ &= (\sqrt{1 - \beta})^2 \mathbf{x}_{t-2} + \sqrt{1 - \beta} \sqrt{\beta} \mathcal{N}(0, I) + \sqrt{\beta} \mathcal{N}(0, I)\end{aligned}$$

## Forward process: a closer look

Let's check if the following  $\alpha$  and  $\beta$  converges:

$$\mathbf{x}_t = \sqrt{1 - \beta} \mathbf{x}_{t-1} + \sqrt{\beta} \mathcal{N}(0, I)$$

$$\begin{aligned}\mathbf{x}_t &= \sqrt{1 - \beta} \mathbf{x}_{t-1} + \sqrt{\beta} \mathcal{N}(0, I) \\ &= \sqrt{1 - \beta} (\sqrt{1 - \beta} \mathbf{x}_{t-2} + \sqrt{\beta} \mathcal{N}(0, I)) + \sqrt{\beta} \mathcal{N}(0, I) \\ &= (\sqrt{1 - \beta})^2 \mathbf{x}_{t-2} + \sqrt{1 - \beta} \sqrt{\beta} \mathcal{N}(0, I) + \sqrt{\beta} \mathcal{N}(0, I) \\ &\dots \\ &= (\sqrt{1 - \beta})^t \mathbf{x}_{t-t} + \dots + (\sqrt{1 - \beta})^2 \sqrt{\beta} \mathcal{N}(0, I) + \sqrt{1 - \beta} \sqrt{\beta} \mathcal{N}(0, I) \\ &+ \sqrt{\beta} \mathcal{N}(0, I)\end{aligned}$$

## Forward process: a closer look

$$\begin{aligned}\mathbf{x}_t = & (\sqrt{1-\beta})^t \mathbf{x}_0 + \cdots + (\sqrt{1-\beta})^2 \sqrt{\beta} \mathcal{N}(0, I) + \sqrt{1-\beta} \sqrt{\beta} \mathcal{N}(0, I) \\ & + \sqrt{\beta} \mathcal{N}(0, I)\end{aligned}$$

- for large  $t$ ,  $(\sqrt{1-\beta})^t$  approaches to 0
- each of  $(\sqrt{1-\beta})^i \sqrt{\beta} \mathcal{N}(0, I)$  is a Gaussian with mean 0 and variance  $(1-\beta)^i \beta$ .
- All these Gaussian can be written as one Gaussian distribution with variance:

$$\sum_{i=0}^t (1-\beta)^i \beta = \beta \frac{1 - (1-\beta)^t}{1 - (1-\beta)} \sim \frac{\beta}{\beta} = 1$$

For large enough  $t$  we converge to the Normal distribution  $\mathcal{N}(0, I)$ .

## Forward process: a closer look

- In practice we do not use a fixed noise  $\beta$
- Linear schedule noise:  $\beta_1 = 0.0001$  and  $\beta_T = 0.02$
- The number of steps is about 10000 (application dependent), but this is slow!
- Recall  $\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I)$ .
- Let's define  $\alpha_t = 1 - \beta_t$  and do the recursive expansion.

## Forward process: a closer look

$$\begin{aligned}\mathbf{x}_t &= \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I) \\&= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \mathcal{N}(0, I) \\&= \sqrt{\alpha_t} (\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \mathcal{N}(0, I)) + \sqrt{1 - \alpha_t} \mathcal{N}(0, I) \\&= (\sqrt{\alpha_t \alpha_{t-1}}) \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t + \alpha_t - \alpha_t \alpha_{t-1}} \sqrt{\beta} \mathcal{N}(0, I) \\&\dots\end{aligned}$$

## Forward process: a closer look

$$\begin{aligned}\mathbf{x}_t &= \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I) \\&= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \mathcal{N}(0, I) \\&= \sqrt{\alpha_t} (\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \mathcal{N}(0, I)) + \sqrt{1 - \alpha_t} \mathcal{N}(0, I) \\&= (\sqrt{\alpha_t \alpha_{t-1}}) \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t + \alpha_t - \alpha_t \alpha_{t-1}} \sqrt{\beta} \mathcal{N}(0, I) \\&\dots \\&= \sqrt{\alpha_t \alpha_{t-1} \cdots \alpha_2 \alpha_1} \mathbf{x}_0 + \sqrt{1 - \alpha_t \alpha_{t-1} \cdots \alpha_2 \alpha_1} \mathcal{N}(0, I) \\&= \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \mathcal{N}(0, I) \\&= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \mathcal{N}(0, I) \\&= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon\end{aligned}$$

## Reverse process: a closer look

Computing the reverse path (reverse denoising distribution) is not, but we know it is diffusion process.

We train a model to approximate the reverse distribution.

Similar to VAEs, the goal is to maximize a lower bound for the likelihood:

$$\log p(\mathbf{x}_0) \geq \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right]$$

A lot of effort goes into simplify this lower bound and exploiting the fact that we are dealing with diffusion process.

## Reverse process: a closer look

At the end of the day, the model should learn the noise added

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0, t, \epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2]$$

What does it mean? How do we actually do training ?

## Reverse process: training

- Sample an image  $\mathbf{x}_0$  and a time step  $t$
- Sample a random noise  $\varepsilon$ .
- Get the noise image at time  $t$ :  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon$
- Feed  $\mathbf{x}_t$  to the neural network.
- The model's predicted noise  $\varepsilon_\theta$  should be close to  $\varepsilon$ .

---

### Algorithm 1 Training

---

- 1: **repeat**
- 2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3:    $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5:   Take gradient descent step on  
     
$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$
- 6: **until** converged

---

## Reverse process: image generation

- Start from a noise image  $\mathbf{x}_T \sim \mathcal{N}(0, I)$ .
- Feed this to the trained neural network to predict a noise  $\epsilon_\theta(\mathbf{x}_T)$
- Given this predicted noise we can do denoising and obtain  $\mathbf{x}_{T-1}$  (we skipped through the math here)
- Repeat the above until we get to  $\mathbf{x}_0$ .

---

### Algorithm 2 Sampling

---

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

---

# Semantic embeddings + diffusion

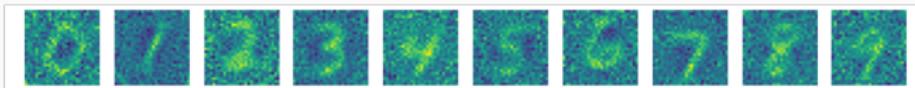
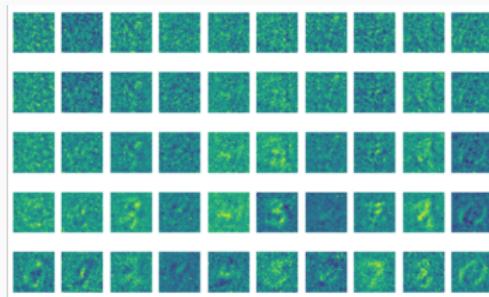
**Text to image synthesis:** Dall-E, Imagen, Stable Diffusion



“A chair that looks like a pineapple”

# Diffusion

A demo for generating digits by training on MNIST.



# Ethical challenges

## How to preserve privacy?

# Generative models and data leakage

Generative models can potentially memorise and regenerate their training data points.

## Extracting Training Data from Diffusion Models

*Nicholas Carlini<sup>\*1</sup> Jamie Hayes<sup>\*2</sup> Milad Nasr<sup>\*1</sup>*

*Matthew Jagielski<sup>+1</sup> Vikash Sehwag<sup>+4</sup> Florian Tramèr<sup>+3</sup>*

*Borja Balle<sup>+2</sup> Daphne Ippolito<sup>†1</sup> Eric Wallace<sup>†5</sup>*

<sup>1</sup>Google   <sup>2</sup>DeepMind   <sup>3</sup>ETHZ   <sup>4</sup>Princeton   <sup>5</sup>UC Berkeley

\*Equal contribution    +Equal contribution    †Equal contribution

# Generative models and data leakage

Generative models can potentially memorise and regenerate their training data points.

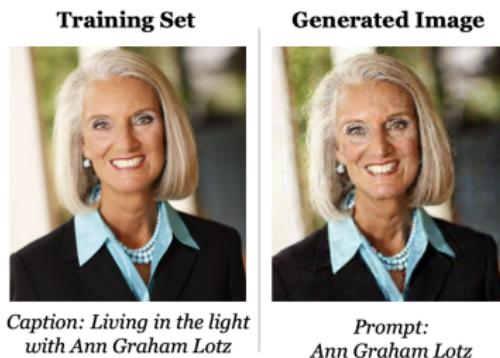


Figure 1: Diffusion models memorize individual training examples and generate them at test time. **Left:** an image from Stable Diffusion’s training set (licensed CC BY-SA 3.0, see [49]). **Right:** a Stable Diffusion generation when prompted with “Ann Graham Lotz”. The reconstruction is nearly identical ( $\ell_2$  distance = 0.031).

## Data leakage

As we saw in the text generation lab, machine learning algorithms are prone to leak information about their training data:

arm towards the viewer. Gregor then turned to look out the window at a barren sister only needed to hear the visitor's first words of greeting and he knew who calm, "I'll get dressed straight away now, pack up my samples and set off. Will again, "seven o'clock, and there's still a fog like this." And he lay there sighing, harder than before, if that was possible, he felt that the lower part of his body a

Here, our generative model revealed entire sentences from the training input. This is a quality issue, but can also be a privacy issue.

## Data leakage

Many modern ML systems trained on user data.

- Smart Compose in Gmail (trained on user emails).
- Generative AI for medical record taking (trained on patient health data).
- Github Copilot trained on public and private repositories.

Even if models do not directly generate private data, it can sometimes be extracted from them.

# Data leakage

Training data extraction attacks can reconstruct verbatim training examples e.g., they can extract secrets such as verbatim social security numbers or passwords.

**Extracting Training Data from Large Language Models**

Nicholas Carlini<sup>1</sup> Florian Tramèr<sup>2</sup> Eric Wallace<sup>3</sup> Matthew Jagielski<sup>4</sup>  
Ariel Herbert-Voss<sup>5,6</sup> Katherine Lee<sup>1</sup> Adam Roberts<sup>1</sup> Tom Brown<sup>5</sup>  
Dawn Song<sup>3</sup> Úlfar Erlingsson<sup>7</sup> Alina Oprea<sup>4</sup> Colin Raffel<sup>1</sup>

<sup>1</sup>*Google* <sup>2</sup>*Stanford* <sup>3</sup>*UC Berkeley* <sup>4</sup>*Northeastern University* <sup>5</sup>*OpenAI* <sup>6</sup>*Harvard* <sup>7</sup>*Apple*

**Abstract**

It has become common to publish large (billion parameter) language models that have been trained on private datasets. This paper demonstrates that in such settings, an adversary can perform a *training data extraction attack* to recover individual training examples by querying the language model.

We demonstrate our attack on GPT-2, a language model trained on scrapes of the public Internet, and are able to extract hundreds of verbatim text sequences from the model's training data. These extracted examples include (public) personally identifiable information (names, phone numbers, and email addresses), IRC conversations, code, and 128-bit UUIDs. Our attack is possible even though each of the above sequences are included in just *one* document in the training data.

We comprehensively evaluate our extraction attack to understand the factors that contribute to its success. Worryingly, we find that larger models are more vulnerable than smaller models. We conclude by drawing lessons and discussing possible safeguards for training large language models.

**1 Introduction**

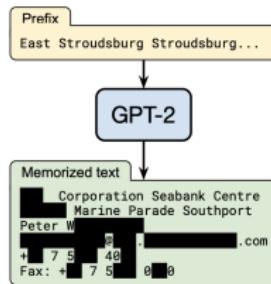


Figure 1: **Our extraction attack.** Given query access to a neural network language model, we extract an individual person's name, email address, phone number, fax number, and physical address. The example in this figure shows information that is all accurate so we redact it to protect privacy.

## The privacy challenge

How do we balance privacy concerns with the desire to train models on as much data as possible?

# Formalizing privacy

There have been many many attempts to formalize what it means for a machine learning algorithm or system to be private.

## Calibrating Noise to Sensitivity in Private Data Analysis

Cynthia Dwork<sup>1</sup>, Frank McSherry<sup>1</sup>, Kobbi Nissim<sup>2</sup>, and Adam Smith<sup>3\*</sup>

<sup>1</sup> Microsoft Research, Silicon Valley. {dwork,mcsherry}@microsoft.com

<sup>2</sup> Ben-Gurion University. kobbi@cs.bgu.ac.il

<sup>3</sup> Weizmann Institute of Science. adam.smith@weizmann.ac.il

Differential Privacy has become the gold standard definition.

Clear theoretical founding, widely used in implemented systems  
(TensorFlow, US Census statistics, Apple User data, etc.)

# Differential privacy

Definition based on notation of neighboring datasets.

**Definition:** A dataset  $\mathbf{X} = [x_1, \dots, x_n]$  is neighbors of a dataset  $\mathbf{X}' = [x'_1, \dots, x'_n]$  if:

$x_i = x'_i$  for all but one value of  $i \in \{1, \dots, n\}$ .

I.e.,  $x_j \neq x'_j$  for a single index  $j$ .

**Alternative but closely related definition:**  $\mathbf{X}$  and  $\mathbf{X}'$  are neighbors if  $\mathbf{X}'$  can be obtained by adding or removing a single data point from  $\mathbf{X}$ .

# Differential privacy

## Definition

An algorithm  $\mathcal{A}$  satisfies  $\epsilon$ -differential privacy if, for any two neighboring datasets  $\mathbf{X}, \mathbf{X}'$ , and any possible output of the algorithm  $\mathbf{z}$ ,

$$\Pr[\mathcal{A}(\mathbf{X}) = \mathbf{z}] \leq e^\epsilon \Pr[\mathcal{A}(\mathbf{X}') = \mathbf{z}].$$

In the context of machine learning,  $\mathcal{A}$  could be the training procedure and  $\mathbf{z}$  could be, e.g., the model weights.

In the context of databases/statistical applications,  $\mathcal{A}$  might implement a simple statistic function like the mean:

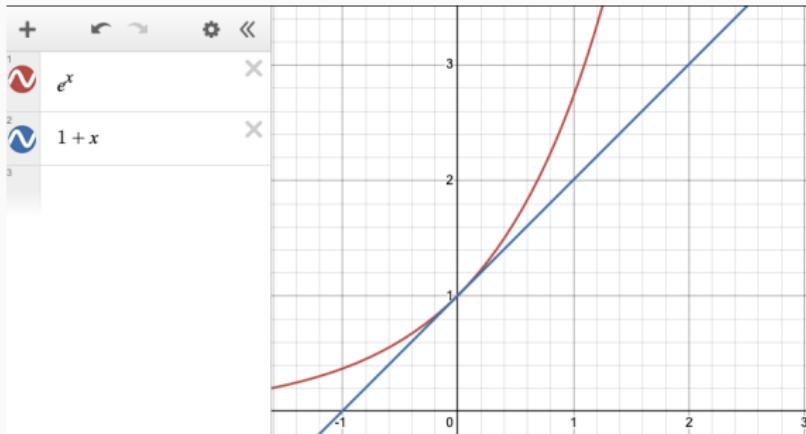
$$\frac{1}{n} \sum_{i=1}^n x_i.$$

# Differential privacy

## Definition

An algorithm  $\mathcal{A}$  satisfies  $\epsilon$ -differential privacy if, for any two neighboring datasets  $\mathbf{X}, \mathbf{X}'$ , and any possible output of the algorithm  $\mathbf{z}$ ,  $\Pr[\mathcal{A}(\mathbf{X}) = \mathbf{z}] \leq e^\epsilon \Pr[\mathcal{A}(\mathbf{X}') = \mathbf{z}]$ .

Think of  $\epsilon$  as a reasonably small constant. E.g.  $\epsilon \in (0, 5]$ . For small  $\epsilon$ ,  $e^\epsilon \approx (1 + \epsilon)$ .



# Differential privacy

## Definition

An algorithm  $\mathcal{A}$  satisfies  $\epsilon$ -differential privacy if, for any two neighboring datasets  $\mathbf{X}, \mathbf{X}'$ , and any possible output of the algorithm  $\mathbf{z}$ ,  $\Pr[\mathcal{A}(\mathbf{X}) = \mathbf{z}] \leq e^\epsilon \Pr[\mathcal{A}(\mathbf{X}') = \mathbf{z}]$ .

In words, differential privacy says that including an individual's data in a dataset  $\mathbf{X}$  can only increase or decrease the probability of observing any particular output by a small factor.

Inherently a property of randomized algorithms. Obtaining differentially private machine learning methods will require **adding randomness to the training process**.

## Differential privacy properties

**Postprocessing property:** If an algorithm  $\mathcal{A}(\mathbf{X})$  is  $\epsilon$ -DP, then  $\mathcal{B}(\mathcal{A}(\mathbf{X}))$  is  $\epsilon$ -DP for any (possibly non-private) algorithm  $\mathcal{B}$ .

**Composition property:** If an algorithm  $\mathcal{A}_1$  is  $\epsilon_1$ -DP and  $\mathcal{A}_2$  is  $\epsilon_2$ -DP, then  $\mathcal{B}(\mathcal{A}_1(\mathbf{X}), \mathcal{A}_2(\mathbf{X}))$  is  $(\epsilon_1 + \epsilon_2)$ -DP.

## Differential privacy

There are many ways to add randomness. Perhaps the most common is noise injection.

**Simple example:** Suppose  $\mathbf{X}$  contains scalar values

$x_1, \dots, x_n \in \{0, 1\}$ . Suppose we want to compute the average,  
$$Q(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n x_i.$$

Naively, adding or removing a point from the dataset changes the average by  $\pm \frac{1}{n}$  with probability 1, so, naively, a mean computation is not differentially private.

## Noise injection

**Differentially Private Estimate of  $Q(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n x_i$ :**

- Generate an appropriate random number  $\eta$ .
- Return  $Q(\mathbf{X}) + \eta$ .

Example =  $\mathbf{X} = \{0, 1, 1, 0, 0, 0\}$ ,  $\mathbf{X}' = \{0, 1, 1, 0, 1, 0\}$ .

Trade-off between privacy and accuracy.

## What type of noise and how much?

### Theorem (Laplace Mechanism)

For a function  $Q$  with sensitivity  $\Delta_Q$ ,

$$\mathcal{A}(\mathbf{X}) = Q(\mathbf{X}) + \text{Lap}(\Delta_Q/\epsilon)$$

is  $\epsilon$ -differentially private.

Sensitivity  $\Delta_Q = \max_{\text{neighboring } \mathbf{x}, \mathbf{x}'} |Q(\mathbf{X}) - Q(\mathbf{X}')|$ .

What is  $\Delta_Q$  for  $Q(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n x_i$ ?

$\text{Lap}(b)$  is a Laplacian random variable with parameter  $b$  (which means variance  $2b^2$ ). PDF is:

$$p_b(\eta) = \frac{1}{2b} e^{-|\eta|/b}$$

# Laplace mechanism analysis

## Theorem (Laplace Mechanism)

For a function  $Q$  with sensitivity  $\Delta_Q$ ,

$\mathcal{A}(\mathbf{X}) = Q(\mathbf{X}) + \text{Lap}(\Delta_Q/\epsilon)$  is  $\epsilon$ -differentially private.

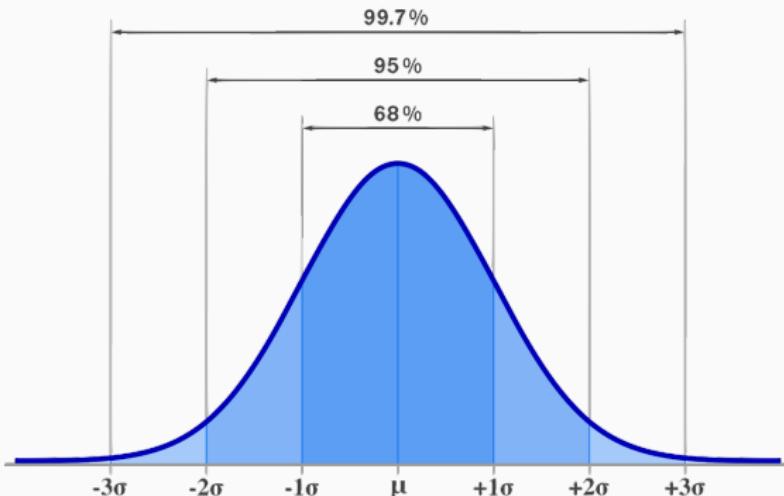
**Proof:** For any possible output  $z$ ,

- $\Pr[\mathcal{A}(\mathbf{X}) = z] = \frac{1}{2(\Delta_Q/\epsilon)} e^{-|Q(\mathbf{X}) - z|/(\Delta_Q/\epsilon)}$
- $\Pr[\mathcal{A}(\mathbf{X}') = z] = \frac{1}{2(\Delta_Q/\epsilon)} e^{-|Q(\mathbf{X}') - z|/(\Delta_Q/\epsilon)}$

$$\begin{aligned}\frac{\Pr[\mathcal{A}(\mathbf{X}) = z]}{\Pr[\mathcal{A}(\mathbf{X}') = z]} &= e^{-(|Q(\mathbf{X}) - z| - |Q(\mathbf{X}') - z|)/(\Delta_Q/\epsilon)} \\ &\leq e^{\frac{|Q(\mathbf{X}) - Q(\mathbf{X}')|}{\Delta_Q/\epsilon}} \leq e^\epsilon.\end{aligned}$$

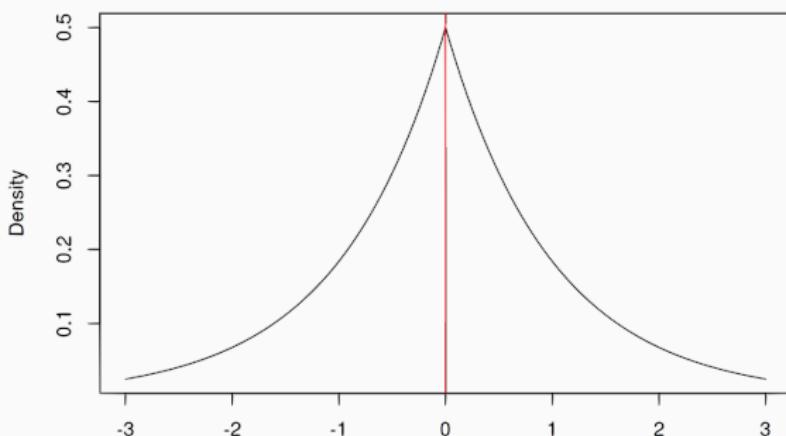
## What do we pay in terms of accuracy?

$Lap(b)$  has standard deviation  $\sqrt{2}b$ . Like Gaussian distribution, Laplace random variables usually fall within a few standard deviations of the mean:



## What do we pay in terms of accuracy?

$Lap(b)$  has standard deviation  $\sqrt{2}b$ . Like Gaussian distribution, Laplace random variables usually fall within a few standard deviations of the mean:



## What do we pay in terms of accuracy?

Standard deviation =  $\sqrt{2} \cdot \frac{\Delta_Q}{\epsilon}$ .

For  $x_1, \dots, x_n \in [0, 1]$ ,  $Q(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n x_i$ , we have that:

$$\Delta_Q = \frac{1}{n}.$$

**Overall error from adding noise:**

$$O\left(\frac{1}{\epsilon n}\right)$$

**Very reasonable if  $n$  is large!**

E.g., if  $n = 10,000$  can get error roughly .001 on mean estimate with privacy parameter  $\epsilon = .1$ .

## What about more complex functions?

In machine learning applications,  $Q$  is an entire training procedure, and the output is vector of parameters.

$$Q(\mathbf{X}, \mathbf{y}) \rightarrow \boldsymbol{\beta} \in \mathbb{R}^d.$$

### Challenges:

- Very hard to estimate the sensitivity to figure out how much noise should be added.
- If some parameters are more sensitive to noise, we could change the models output drastically.

## Differentially private (stochastic) gradient descent

**Main idea:** Typically  $Q(\mathbf{X}, \mathbf{y})$  is computed by running gradient descent on a loss function  $L(\beta)$ . Instead of directly adding noise to  $Q(\mathbf{X}, \mathbf{y})$ , add noise at each step of gradient descent.

### Basic Gradient descent algorithm:

- Choose starting point  $\beta^{(0)}$ .
- For  $i = 0, \dots, T$ :
  - $\beta^{(i+1)} = \beta^{(i)} - \eta \nabla L(\beta^{(i)})$
- Return  $\beta^{(T)}$ .

## Differentially private (stochastic) gradient descent

Typical loss function in machine learning have finite sum structure.

$$L(\beta) = \sum_{j=1}^n \ell(\beta, \mathbf{x}_j, y_j)$$

By linearity:

$$\nabla L(\beta) = \sum_{j=1}^n \nabla \ell(\beta, \mathbf{x}_j, y_j)$$

Looks just like our mean estimation problem! Can bound the contribution of each data example  $(\mathbf{x}_j, y_j)$  to the gradient to get a sensitivity, then add noise.

## Differentially private (stochastic) gradient descent

Due to a 2016 paper by Martín Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, Li Zhang.

### DP-SGD:

- Choose starting point  $\beta^{(0)}$ .
- For  $i = 0, \dots, T$ :
  - $\beta^{(i+1)} = \beta^{(i)} - \eta(\nabla L(\beta^{(i)}) + \mathbf{r}_i)$
- Return  $\beta^{(T)}$ .

Above each  $\mathbf{r}_i$  is a random Gaussian vector.

**Leading way to incorporate privacy into training machine learning models. Implemented natively, e.g., in TensorFlow.**