# Deep Generative Models Chapter 2: Data Generation Problem

#### Ali Bereyhi

#### ali.bereyhi@utoronto.ca

Department of Electrical and Computer Engineering University of Toronto

Summer 2025

### **Beyond Text**

In Chapter 1 . . .

We considered only text generation which is one possible modality

A data could though come from different natures  $\equiv$  modalities

- Text data typically has a form of linguistic modality
  - → We can treat it as a sequence of unit objects
- Images incorporate their information in their visual configuration
  - → We represent them as arrays of pixels
- In audio data information is incorporated temporally as well
- Video encodes information in both temporal and visual form

• • • •

### Beyond Text: Generic Datatype

At the end of day, each sample is presented in a mathematical form

- We modeled text as a sequence tokens
  - $\downarrow$  It is a sequence of integers  $\equiv$  token indices
- Image are multi-channel tensors of pixel values
  - $\downarrow$  A CIFAR-10 image is a  $3 \times 32 \times 32$  tensor of 8-bit pixels
- Audio files are sequences of framed samples
  - → We sample with Nyquist rate and put them in fixed-size frames
- Videos are sequences of tensor frames
  - L→ Each frame contains multiple image tensors

#### Generic Datatype

We can denote a generic data sample by a mathematical object  $\boldsymbol{x}$ 

### Data Space

#### Data Space

Data sample x comes from a known data space X, i.e.,  $x \in X$ 

- We have a text  $x \equiv x_1, \ldots, x_T$  from a vocabulary of size I
  - $\vdash$  Each token  $x_t \in \{0, \ldots, I-1\}$
  - $\downarrow$  Data samples belong to data space  $\mathbb{X} = \{0, \dots, I-1\}^T$
- An  $n \times n$  pixel RGB image contains three color-maps
  - ${\scriptstyle {\rm L}}{\scriptstyle {\rm >}}$  Each color map contains  $n \times n$  pixels
  - $\vdash$  Each pixel is represented by 8 bits, i.e.,  $\{0, \ldots, 255\}$
  - $\downarrow$  Data sample belongs to data space  $X = \{0, \dots, 255\}^{3 \times n \times n}$

Our data is a  $5 \times 5$  pixel digit

• Each pixel is either  $0 \equiv \text{off}$  or  $1 \equiv \text{on}$ 

We can represent each sample by a matrix x coming from  $\mathbb{X} \{0, 1\}^{5 \times 5}$ 



Our data is a  $5 \times 5$  pixel digit

• Each pixel is either  $0 \equiv \text{off}$  or  $1 \equiv \text{on}$ 

We can represent each sample by a matrix x coming from  $\mathbb{X}\left\{0,1\right\}^{5\times 5}$ 



Our data is a  $5 \times 5$  pixel digit

• Each pixel is either  $0 \equiv \text{off}$  or  $1 \equiv \text{on}$ 

We can represent each sample by a matrix x coming from  $\mathbb{X}\left\{0,1\right\}^{5\times 5}$ 



Our data is a  $5 \times 5$  pixel digit

• Each pixel is either  $0 \equiv \text{off}$  or  $1 \equiv \text{on}$ 

We can represent each sample by a matrix x coming from  $\mathbb{X} \{0, 1\}^{5 \times 5}$ 



### Valid Data Points

Though data space represents where data comes from ....

not all points in the data space are valid data samples







### Valid Data Points $\equiv$ Data Distribution

- + How can we then specify valid data points in the data space?
- We can define a distribution for data

### Data Distribution

Data distribution P(x) is the distribution by which data samples are drawn from data space X

- + How does it help defining valid data points?
- Well! The probability of each sample can define its validity
  - → Invalid data points happen with probability zero
  - → Frequent data points happen with higher probability
  - → Non-frequent data points happen with less probability

### Data Distribution vs Dataset

Data distribution describes the statistical viewpoint on data

In this viewpoint: each sample that we see

- → is drawn randomly from data distribution
- ↓ this is indeed why we call it a data sample

Dataset

A collection of samples drawn from data distribution P(x)

 $\sim$ 

## Example: $5 \times 5$ Pixel Digits

For  $5\times 5$  pixel digits

#### └→ invalid matrices happen with probability zero



$$\Rightarrow P(x = \begin{bmatrix} 1 & 1 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \end{bmatrix}) = 0$$

For  $5\times 5$  pixel digits: digits might be given at the same chances



For  $5\times 5$  pixel digits: some digits might be given at higher chances

→ valid matrices happen with different non-zero probabilities



# Data Distribution: Complexity

- + Wouldn't it be easier to work with the set of valid data-points instead of defining data distribution?
- Absolutely No!

Data distribution enables us to capture

- the pattern hidden in the data
  - └→ E.g., different samples can occur wit different frequencies
- a feasible way to model intractable data complexity
  - → Just think about a realistic dataset like CIFAR-10

### Example: CIFAR-10

In CIFAR-10, we have  $32 \times 32$  RGB images with uint8 pixels labeled by

airplanes  $\equiv 0, \ldots, trucks \equiv 9$ 

A data sample is hence represented by

 $(x, y) \in \mathbb{X} = \{0, \dots, 255\}^{3 \times 32 \times 32} \times \{0, \dots, 9\}$ 

+ How many of such samples we have?

# possible points in  $X = 256^{3072} \cdot 10 = 2^{24576} \cdot 10 > 10^{2458}$ 

We could have had all samples by now if we would have kept collecting

since big-bang with a rate more than  $10^{2440}$  samples/sec!

**Deep Generative Models** 

#### Data Distribution

# Data Distribution is Everything!

What we have in CIFAR-10 dataset is not even a drop of the ocean!

Defining data distribution helps us think about it rigorously, e.g., in CIFAR-10



Distribution is All We Need

Indeed, learning data distribution is the whole thing we do in machine learning!

Using this notion, whatever we've done is formulated in universal framework

We study this universal framework in the next section

### Problem of Data Generation: Most Generic Form

+ Can we present the problem of data generation formally then?

- Yes!

### Data Generation

We look for a generation mechanism which can learn sampling from a data distribution after observing a limited set of samples drawn from that distribution



We will realize this mechanism with different tricks throughout the semester