

Generative Deep Learning

Naeemullah Khan

naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله
للغات والتكنولوجيا

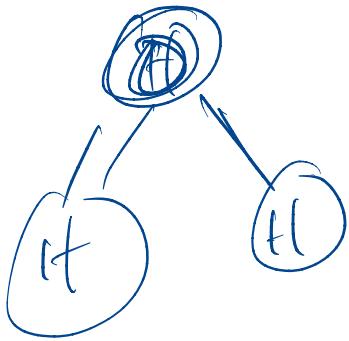
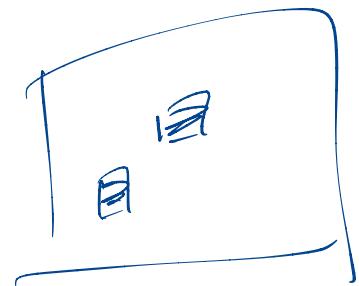
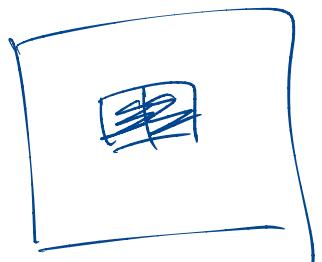
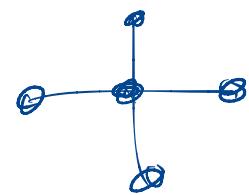
King Abdullah University of
Science and Technology



LMH
Lady Margaret Hall

KAUST Academy
King Abdullah University of Science and Technology

July 11, 2023



① Deep Unsupervised Generative Model

② Reinforcement Learning

③ Graph Neural Network

- ① Classical ML (Regression, logistic Regress)
- ② NN /
- ③ CNN /
- ④ Back propagation
- ⑤ Calculus / linear Algebra / optimised
- ⑥ PyTorch / TensorFlow

ML

Supervised
Learning

- ① Segmentation
- ② Structure Data

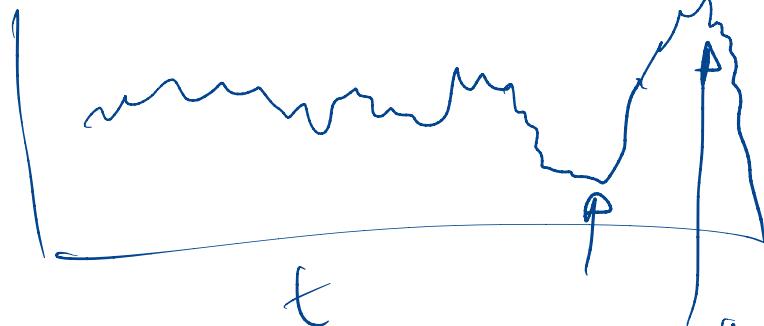
x, y

$$y = f_0(x)$$

$$x \in \mathbb{R}^n$$

\hat{f}

Supervised Learning

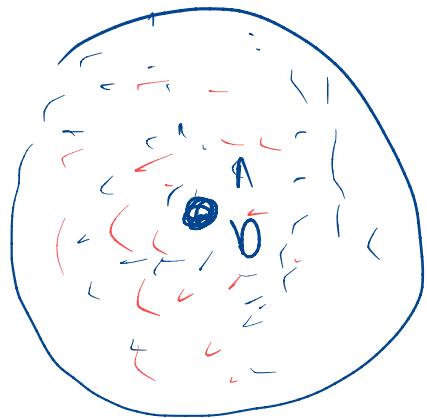


$$x \in \mathbb{R}^d$$

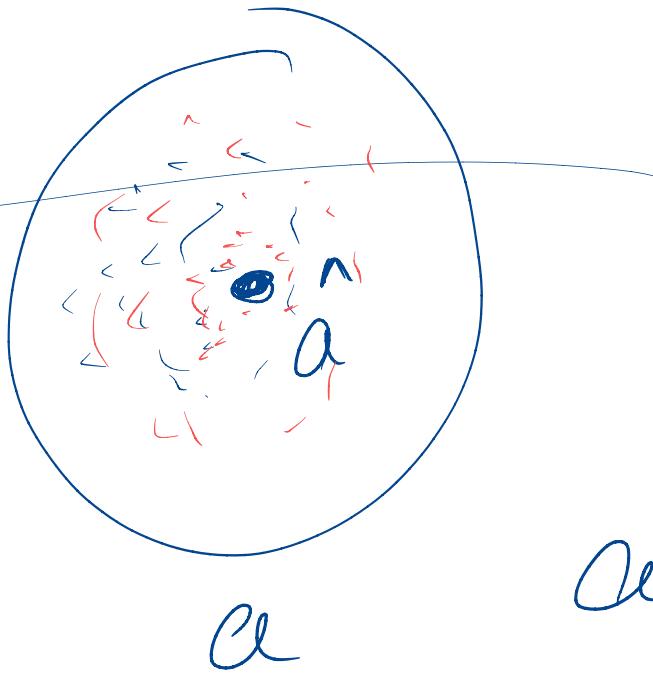
d is large

$$x \not\in \mathcal{D}_X$$

acc



B



a

3×10^7 bysl

~~3KB~~ 30 KB

$$x_i \in \mathbb{R}^3$$

$$x_i \in a$$

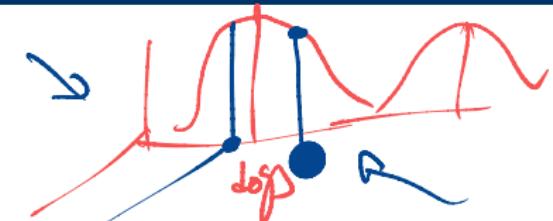
EGNC(0;1)

$$\alpha_i z \hat{a} + \theta_i$$

$$b_j = \hat{b} + \epsilon_j$$

3×2 byte

- ① comprehension
- ② structure → supervised tasks
- ③ Generation ✓



- ▶ Can you imagine an image of dog?
- ▶ Now, can you imagine an image of that dog driving a car?
- ▶ How are we able to this?
- ▶ We excel at extracting knowledge from the data we observe and perform complex reasoning based on it



$$(2^{56})^{16}$$

$$\leftarrow$$

$$x_1 \dots x_{16} \in \{0, 1\}^{56}$$

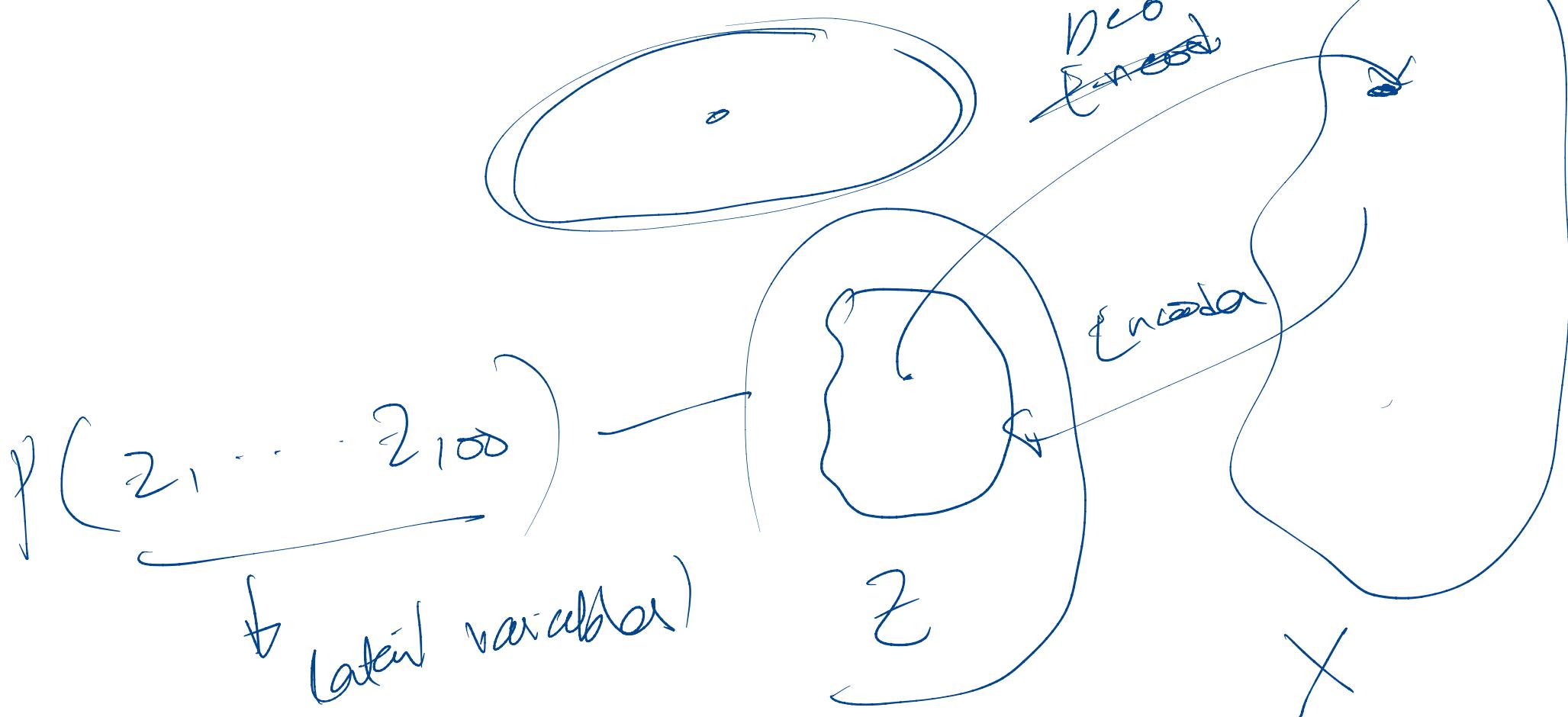
$$X \in \{0, 1\}^{56}$$

$$x_1, x_2 \in \{0, 1\}^{56}$$

$$x_1 \dots x_{16} \in \{0, 1\}^{56}$$

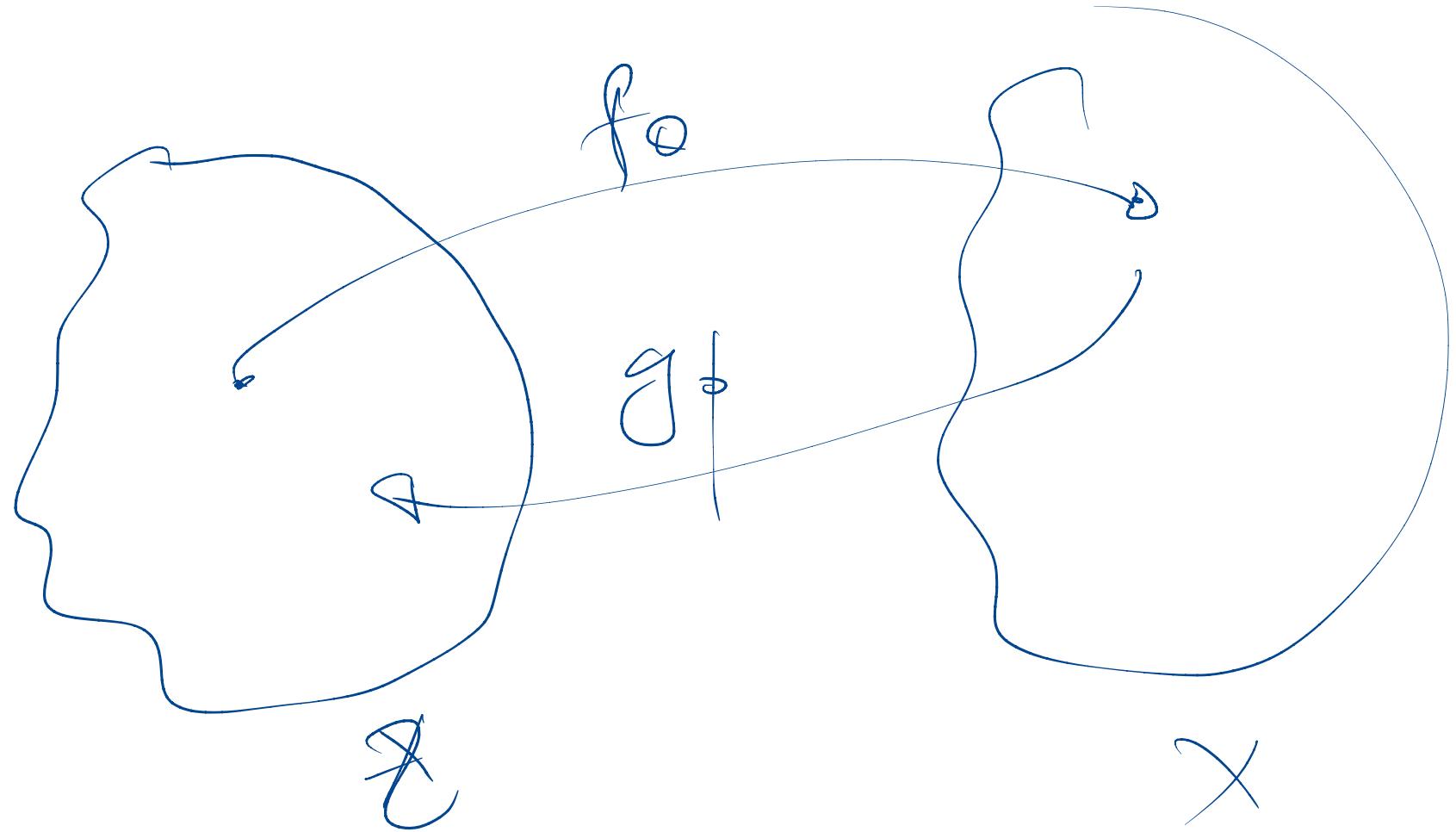
$$P(x_1, \dots, x_{10^6}) = \prod_{i=1}^{10^6} P(x_i) \times$$

P



$$P(z_1, \dots, z_{10^6}) = \prod_{i=1}^{10^6} p(z_i)$$

Normalized flow



Auto Regressive Models

$$P(x_1, x_2, x_3, \dots, x_6) = P$$

$$P(x_1) P(x_2 | x_1) \cdot P(x_3 | x_2, x_1) \cdot P(x_4 | x_3, x_2, x_1)$$

$$P(x_5 | x_4, x_3, \dots, x_1) \cdot P(x_6 | x_5, \dots, x_1)$$

$$P(x_1, x_2) \cdot P(x_3 | x_2, x_1) = P(x_1, x_2, x_3)$$

$$P(x_1, x_2, x_3, x_4), P(x_5, x_6, x_7, x_8)$$

$$P(x_1, x_2) = P(x_2 | x_1) \cdot P(x_1)$$

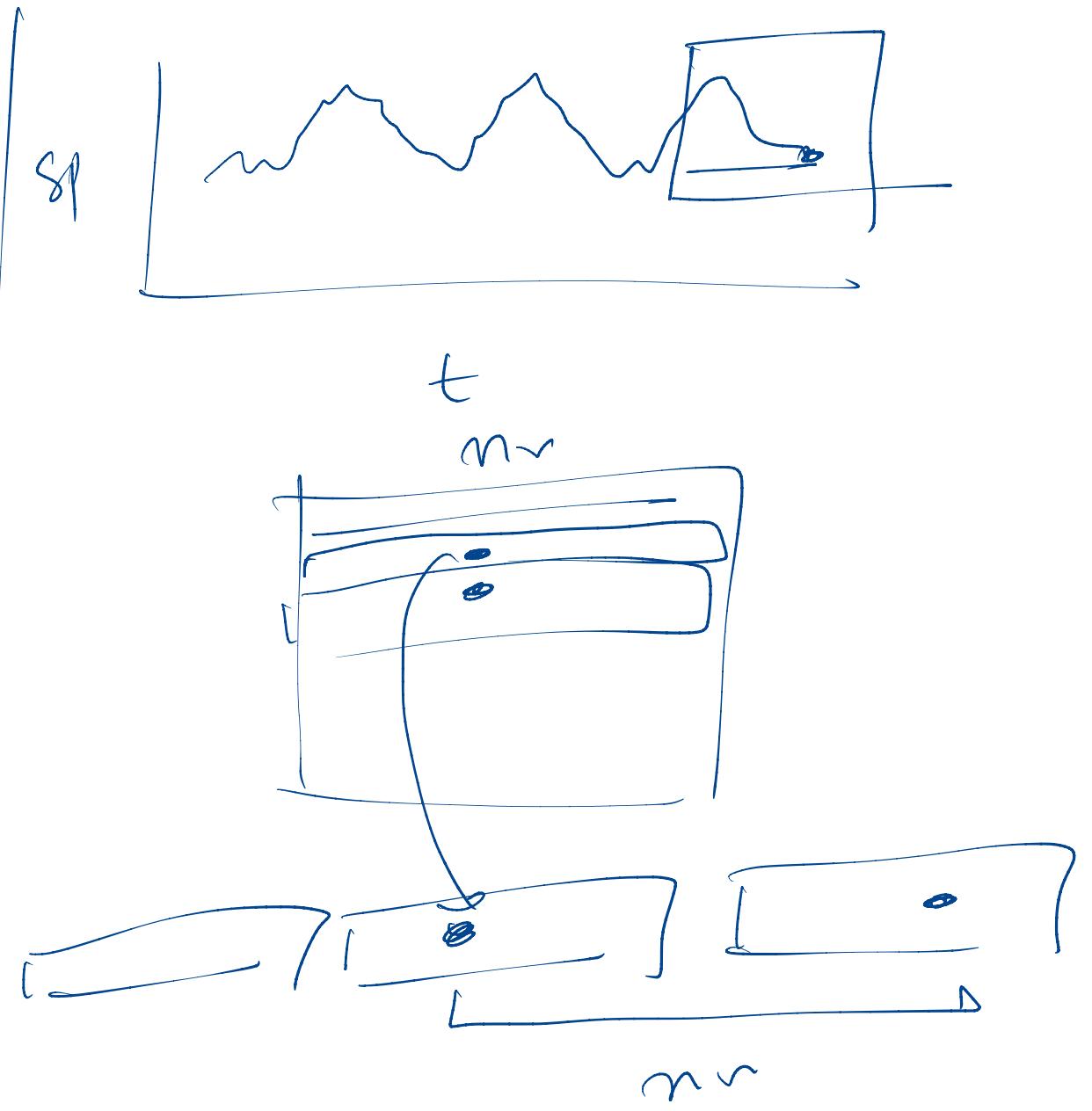
$$P(x_1 | x_1) = \frac{P(x_1, x_2)}{P(x_2)}$$

Bayes Rule

$$P(x_1, \dots, x_M) \\ = \prod_{i=1}^M P(x_i | \underbrace{x_{i+1}, \dots, x_M}_{\text{fixed}})$$

with some except

$$\prod_{i=1}^M P(x_i | h(x_i))$$



$\{x_1, x_2, \dots, x_m\}$

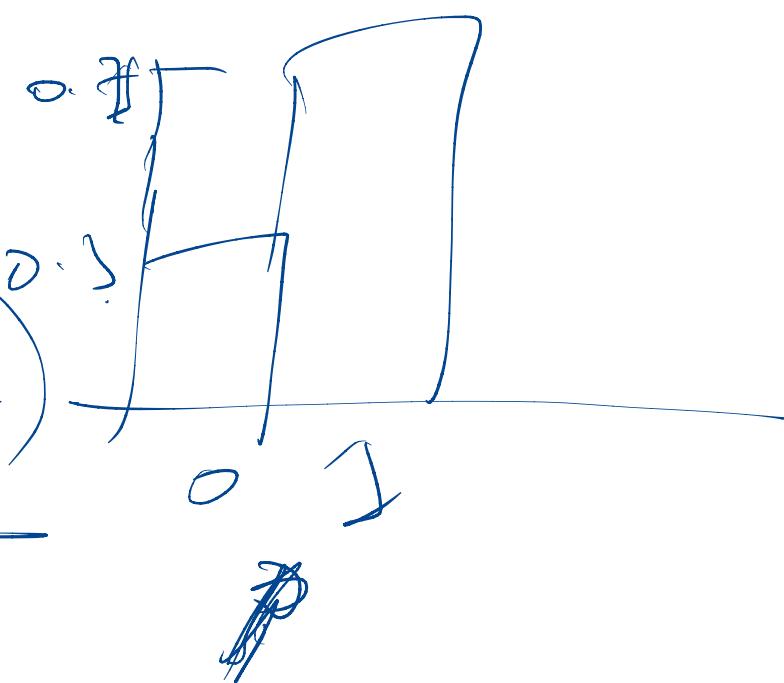
$x \rightarrow y$

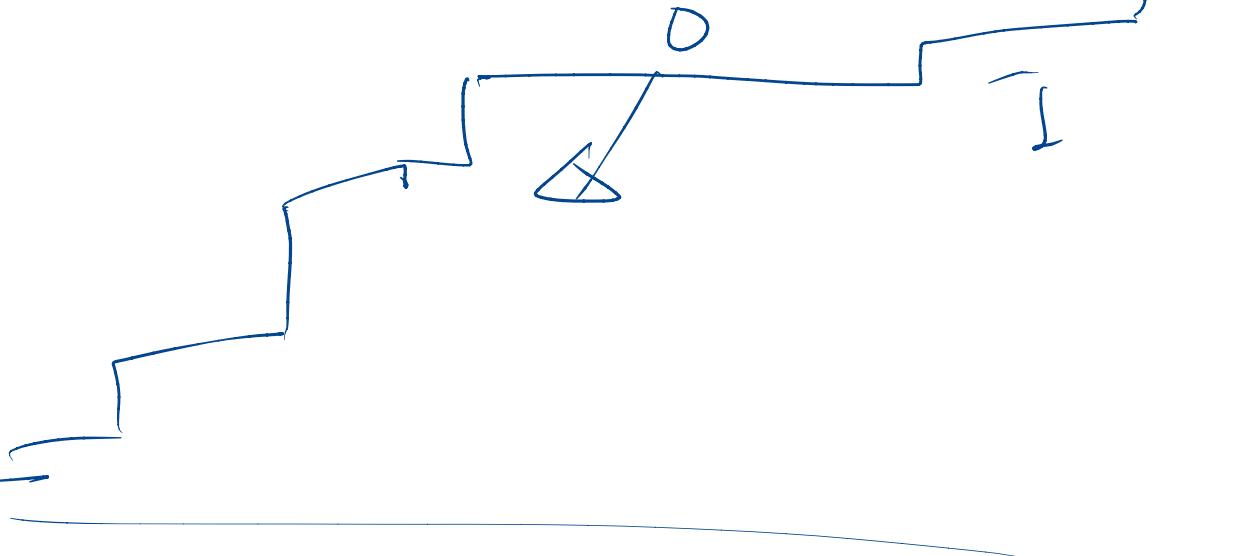
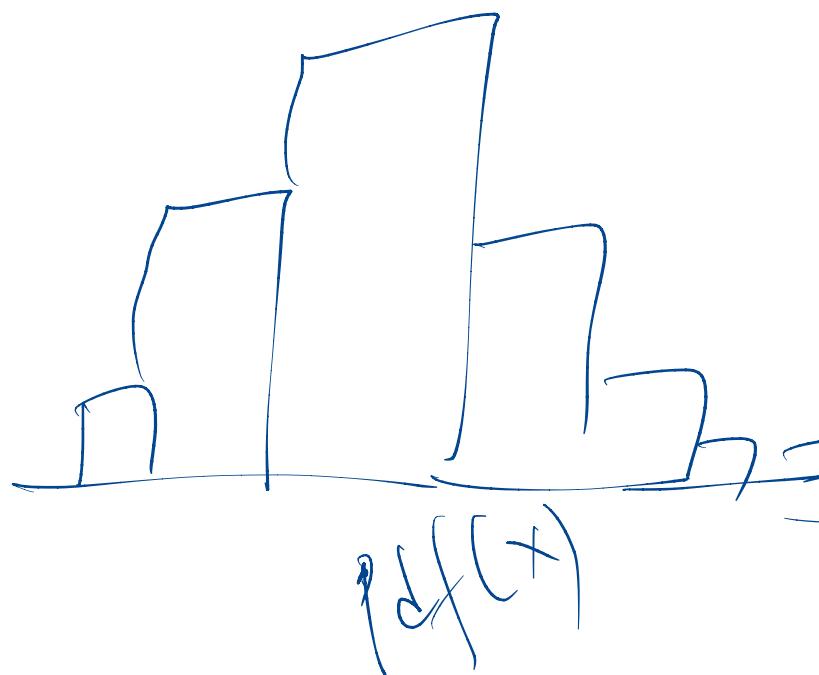
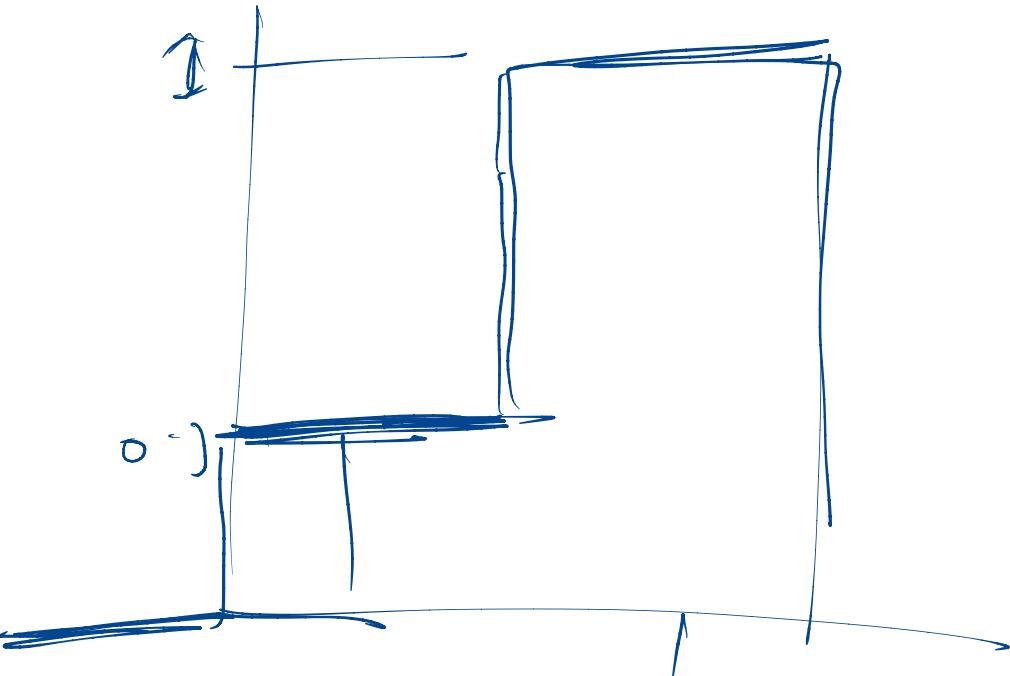
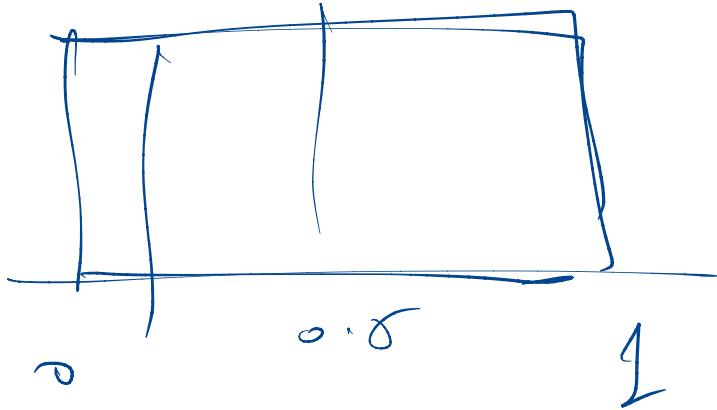
$$P(x_1) = P(x_2 | x_1) \cdot P(x_3 | \underline{x_2}) \cdot P(x_4 | \underline{x_2, x_3})$$

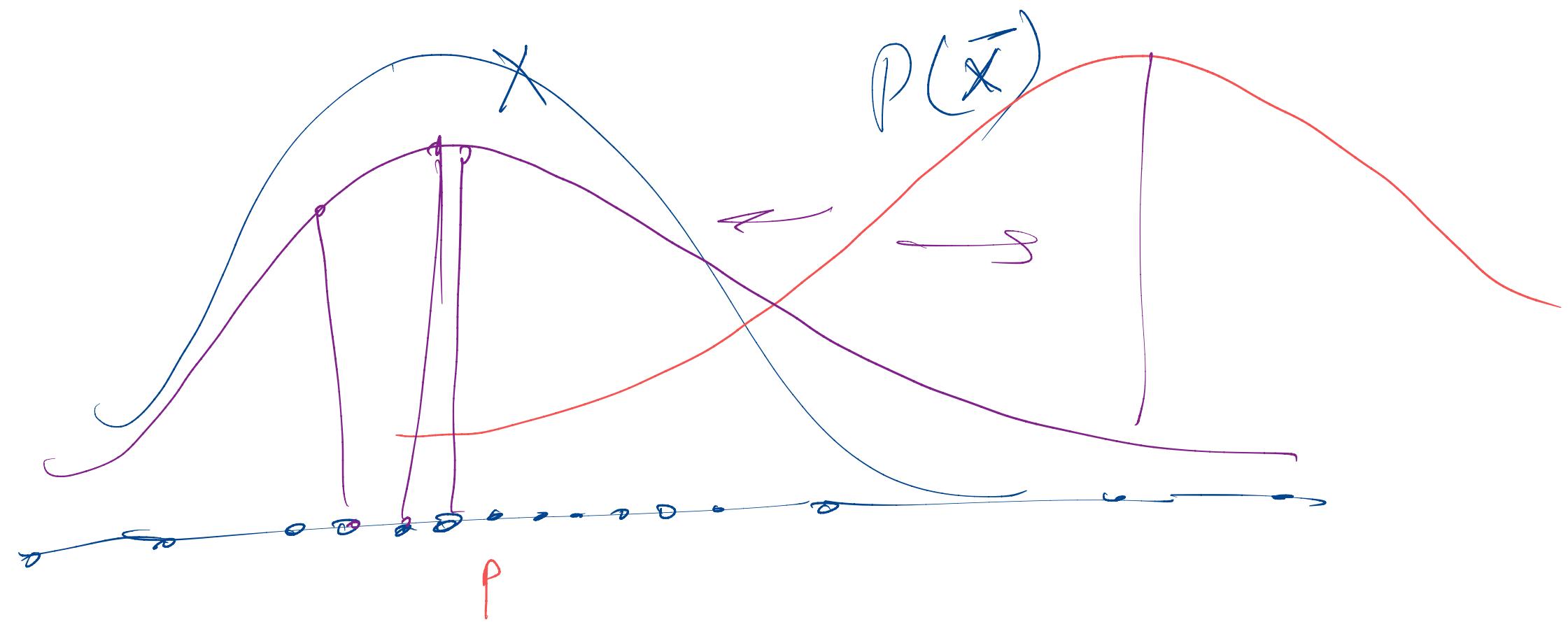
$$P(x_2 | \cdot) \sim \mathcal{N}(\alpha_0 + \alpha_1 x_1) \text{ if } f$$

$$P(x_3 | x_1, x_2) \sim \mathcal{N}(\beta_0 + \beta_1 x_1 + \beta_2 x_2)$$

#







How can we generate such images/data through AI

What we would need? Let's think about it step by step.

- ▶ Some prior data
- ▶ A model which will "learn" the data
- ▶ A method through which the model will learn
- ▶ A method to generate new data from the trained model

This what generative deep learning is all about and what we will be covering in this course.

$$P_2 \frac{1}{\sqrt{2\pi}} \exp^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\underline{x(1)}, \underline{x(2)}, \dots, \underline{x(n)}$$

$$P(x_1, x_2, \dots, x_N) = \prod_{i=1}^N P(x_i)$$

iif

\Rightarrow likelihood

$$\max_u \prod_{i=1}^N P(x_i)$$

$$\max_a \sum_{i=1}^N \log(P_{X_i})$$

$$\max_u -N \left(\log \frac{1}{\sqrt{2\pi}} \right) + \sum_{i=1}^N \frac{(x_i - u)^2}{2}$$

LL \rightarrow sigmoid

NLL \rightarrow avg min u

$$\log \frac{1}{\sqrt{2\pi}} \exp \frac{-(x-u)^2}{2}$$

$$\sum_{i=1}^N (x_i - u)^2$$

$$\begin{aligned}
 & \frac{1}{\sqrt{2\pi} \sigma} \exp \left(-\frac{(x-u)^2}{2\sigma^2} \right) \\
 & \quad \left. \frac{d}{du} \sum_{i=1}^n \frac{(x_i - u)^2}{2} \right|_{u=\bar{x}} \\
 & \quad \sum_{i=1}^n \frac{\partial}{\partial u} \frac{(x_i - u)^2}{2} \\
 & \quad \sum_{i=1}^n (x_i - u) = 0 \\
 & \quad \sum_{i=1}^n x_i = n \bar{x} \Rightarrow \bar{x} \sum_{i=1}^n 1 = n \bar{x} \\
 & \boxed{u = \frac{1}{N} \sum_{i=1}^n x_i}
 \end{aligned}$$

Introduction (cont.)

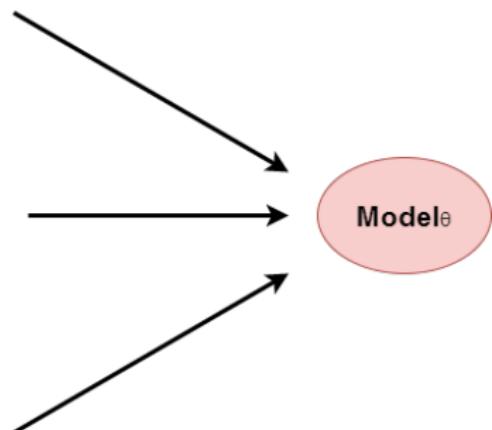
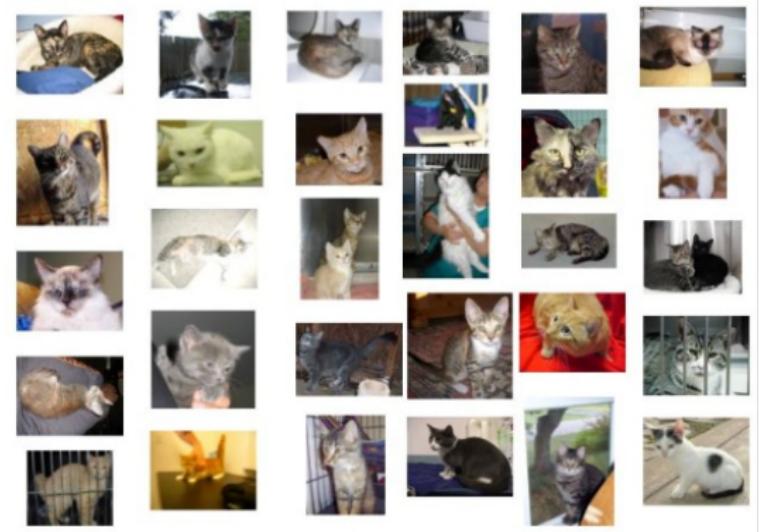


Figure 2: The model learns the cat images data

Introduction (cont.)

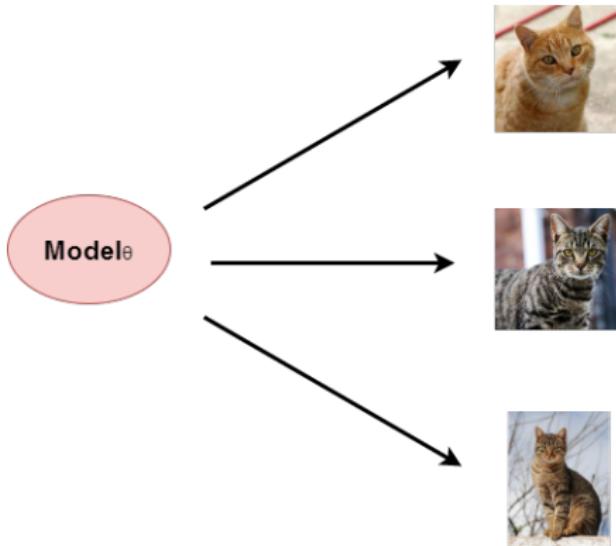


Figure 3: The trained model generates new cat images

Statistical generative models are generative models that are learned from data.

- ▶ Data $\mathcal{X} \sim \mathcal{P}(\mathcal{X})$
- ▶ Model $\theta \sim \mathcal{P}(\theta)$
- ▶ Optimize model s.t. $\mathcal{P}(\theta) \sim \mathcal{P}(\mathcal{X})$
- ▶ Generate $\mathcal{X}_{new} \sim \mathcal{P}(\theta)$

Discriminative Models

- ▶ E.g.: Classify between cat and dog images
- ▶ Doesn't really need to learn data distribution $\mathcal{P}(\mathcal{X})$
- ▶ Goal is to learn $\mathcal{P}(\mathcal{Y} = \text{dog} | \mathcal{X} = x)$

Generative Models

- ▶ E.g.: Generate new cat images
- ▶ Need to learn data distribution $\mathcal{P}(\mathcal{X})$
- ▶ Goal is to maximize $\mathcal{P}(\mathcal{X} = x_{\text{new}})$

Progress on Face Generation



Figure 4: Progress in face generation over the years¹

¹<https://twitter.com/tamaybes/status/1450873331054383104>

Completing Incomplete Content

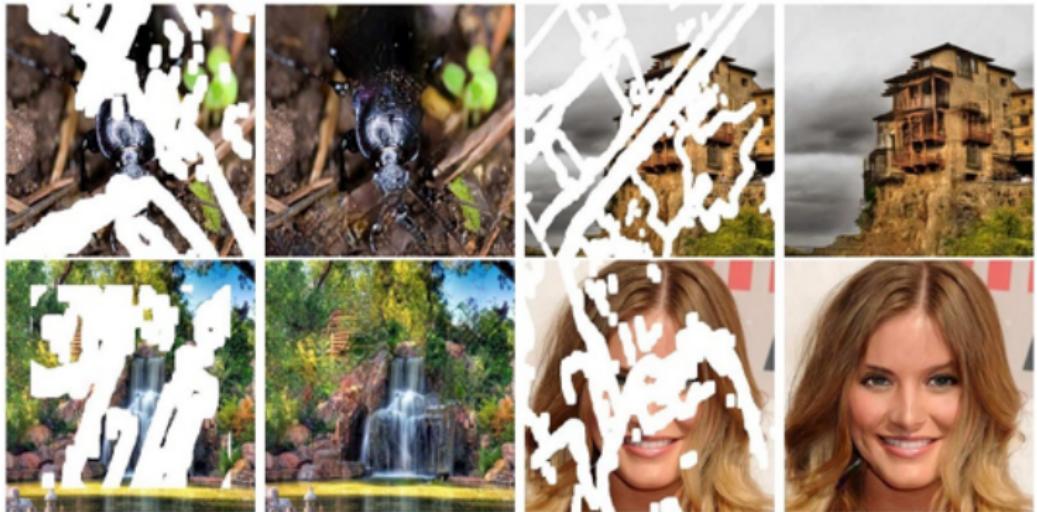


Figure 5: Complete missing patches in image

Completing Incomplete Content (cont.)



Figure 6: Improve image quality

Completing Incomplete Content (cont.)



Antic, 2020

Figure 7: Color b/w images

Completing Incomplete Content (cont.)

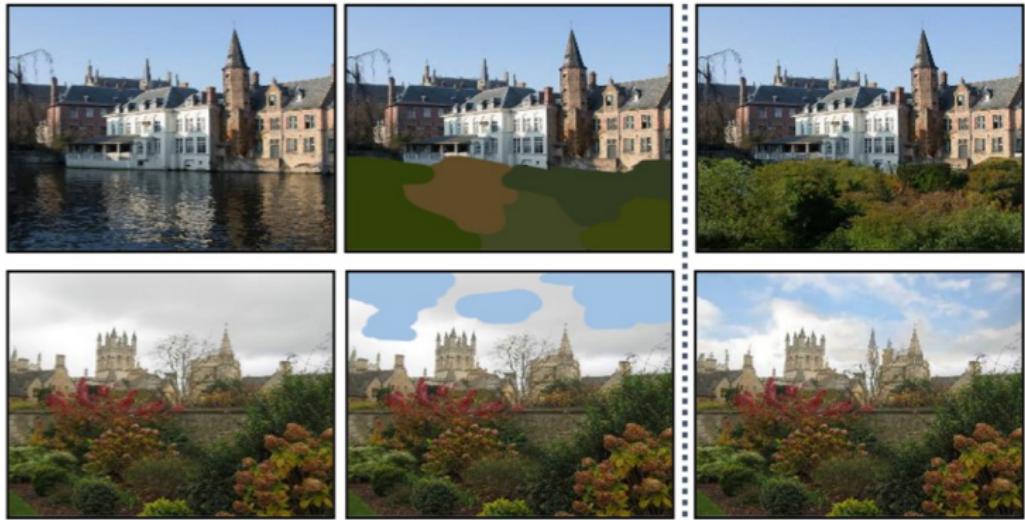


Figure 8: Editing based on strokes in image

Completing Incomplete Content (cont.)

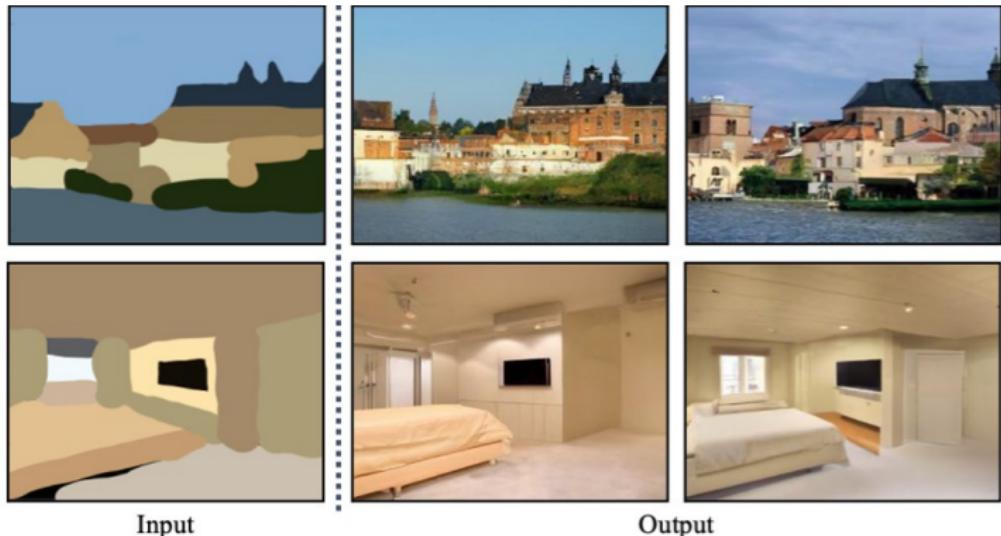


Figure 9: Converting strokes to proper images

Text Generation

InferKit DEMO

Generate Options

Recent advances in generating new content through artificial intelligence |

Learn more in [the docs](#).

Length to generate ⓘ 560

Try to include these words ⓘ

Type some words

Start at beginning ⓘ

[Advanced Settings »](#)

Generate Text

X

File

Figure 10: Completing Text. <https://app.inferkit.com/demo>

What is Deep Unsupervised Learning?

- Capturing rich patterns in raw data with deep networks in a **label-free** way
 - Generative Models: recreate raw data distribution
 - Self-supervised Learning: “puzzle” tasks that require semantic understanding
- But why do we care?



Geoffrey Hinton

(in his 2014 AMA on Reddit)

“The brain has about 10^{14} synapses and we only live for about 10^9 seconds. So we have a lot more parameters than data. This motivates the idea that we must do a lot of unsupervised learning since the perceptual input (including proprioception) is the only place we can get 10^5 dimensions of constraint per second.”



Yann LeCun

Need tremendous amount of information to build machines that have common sense and generalize

[LeCun-20161205-NeurIPS-keynote]

■ “Pure” Reinforcement Learning (**cherry**)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

LeCake

■ Supervised Learning (**icing**)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

■ Unsupervised/Predictive Learning (**cake**)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**



■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)

“Ideal Intelligence”

“Ideal Intelligence” is all about compression (finding all patterns)

- Finding all patterns = short description of raw data (low Kolmogorov Complexity)
- Shortest code-length = optimal inference (Solomonoff Induction)
- Extensible to optimal action making agents (AIXI)

Aside from theoretical interests

- Deep Unsupervised Learning has many powerful applications
 - Generate novel data
 - Conditional Synthesis Technology (WaveNet, GAN-pix2pix)
 - Compression
 - Improve any downstream task with un(self)supervised pre-training
 - Production level impact: Google Search powered by BERT
 - Flexible building blocks

Generate Images



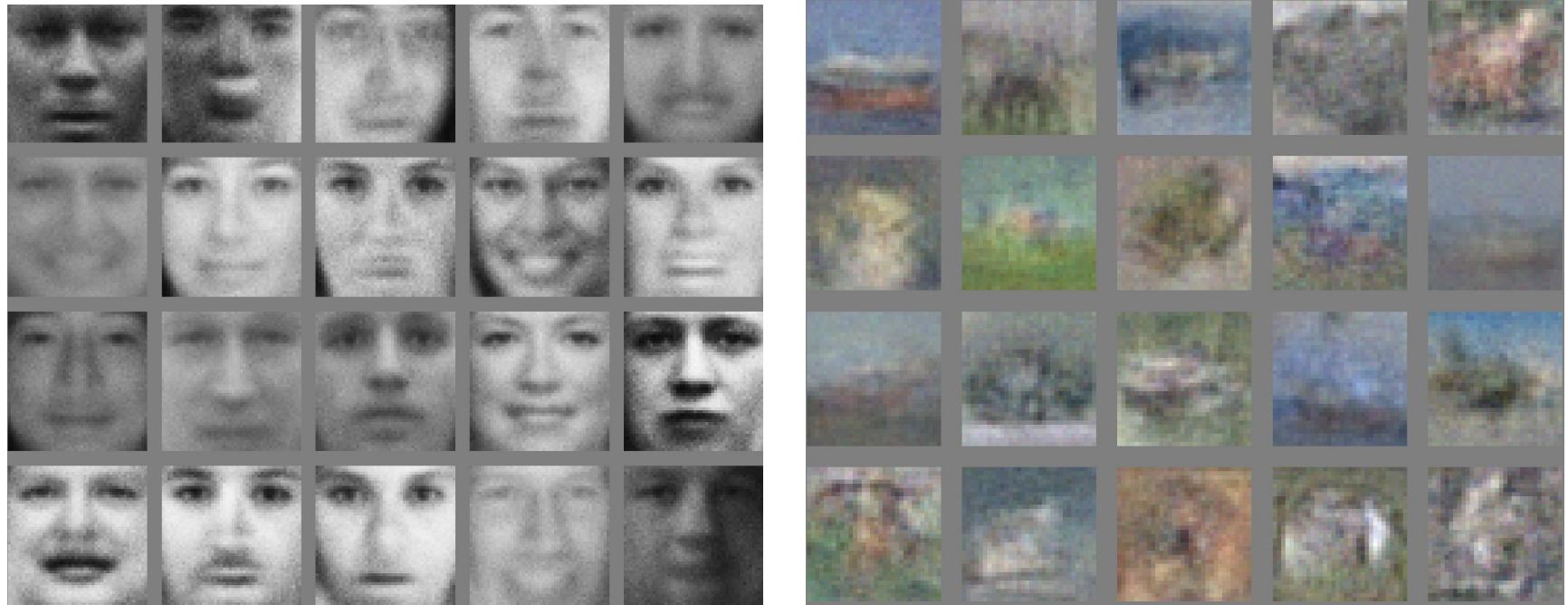
[Deep Belief Nets, Hinton, Osindero, Teh, 2006]

Generate Images



[VAE, Kingma and Welling, 2013]

Generate Images



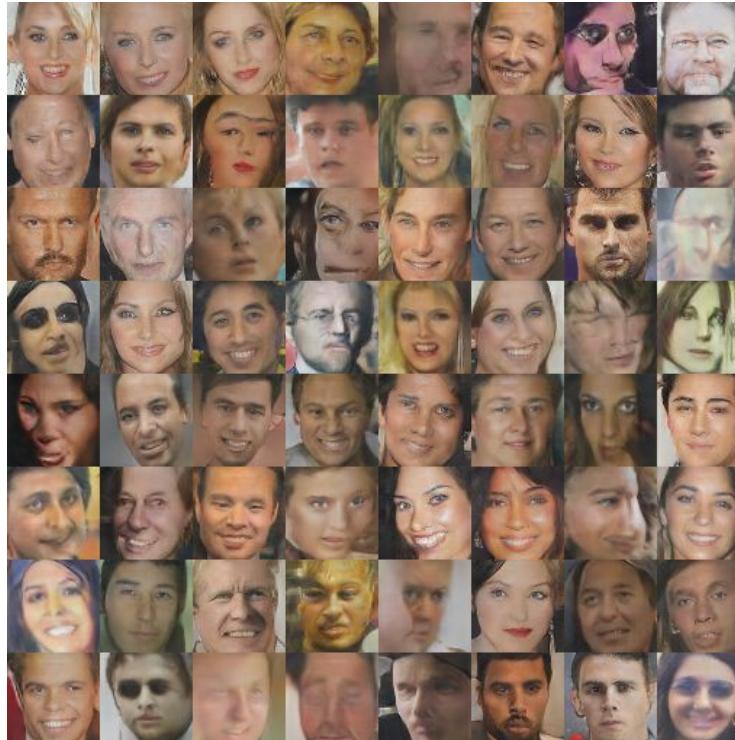
[GAN, Goodfellow et al. 2014]

Generate Images



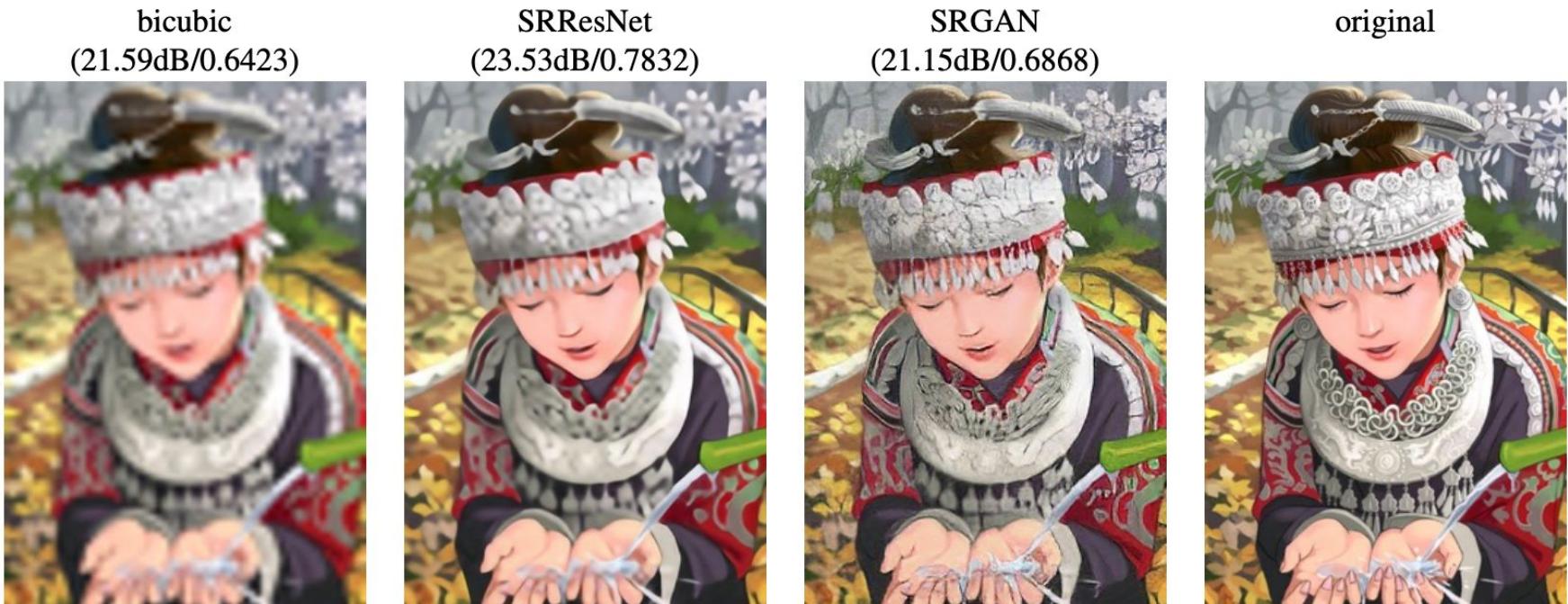
[DCGAN, Radford, Metz, Chintala 2015]

Generate Images



[DCGAN, Radford, Metz, Chintala 2015]

Generate Images



[Ledig, Theis, Huszar et al, 2017]

Generate Images



[CycleGAN: Zhu, Park, Isola & Efros, 2017]

Generate Images



[BigGAN, Brock, Donahue, Simonyan, 2018]

Generate Images



[StyleGAN, Karras, Laine, Aila, 2018]

Generate Audio



1 Second

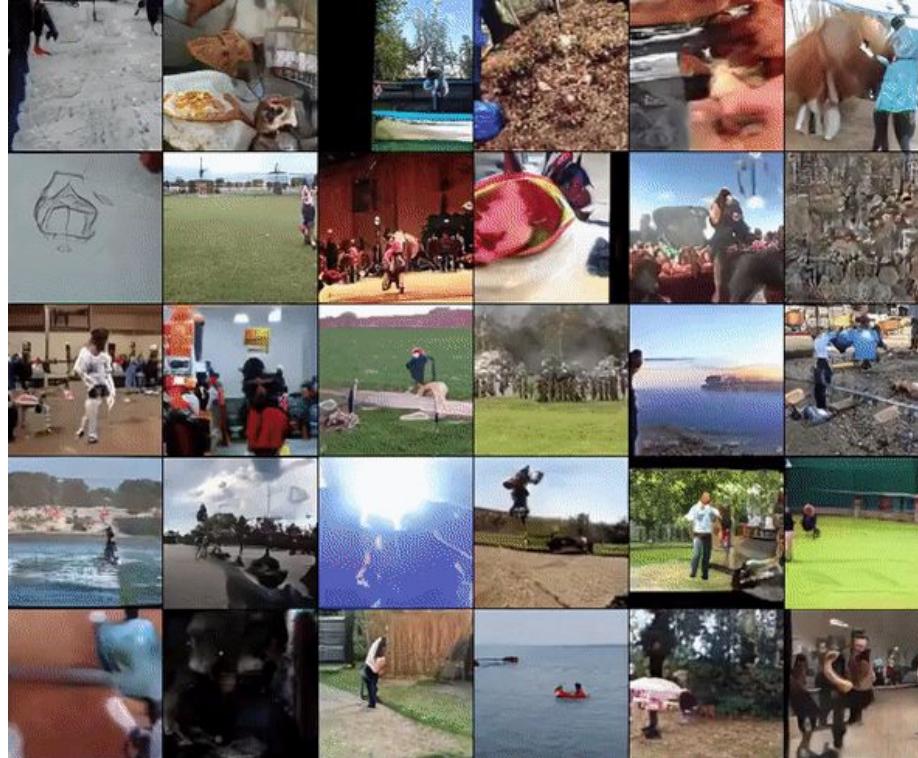


Parametric

WaveNet

[WaveNet, Oord et al., 2018]

Generate Video



DVD-GAN: Adversarial Video Generation on Complex Datasets, Clark, Donahue, Simonyan, 2019

Generate Text

PANDARUS:

Alas, I think he shall be come approached and the day
When little strain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

[Char-rnn, karpathy, 2015]

Generate Math

```
\begin{proof}
```

We may assume that \mathcal{F} is an abelian sheaf on \mathcal{C} .

Given a morphism $\Delta : \mathcal{F} \rightarrow \mathcal{F}$ to \mathcal{F} is injective and let \mathfrak{q} be an abelian sheaf on X .

Let \mathcal{F} be a fibered complex. Let \mathcal{F} be a category.

```
\begin{enumerate}
```

Given $\text{set in construction phantom}$ {Lemma}

label{lemma-characterize-quasi-finite}

Let \mathcal{F} be an abelian quasi-coherent sheaf on \mathcal{C} .

Let \mathcal{F} be a coherent \mathcal{O}_X -module.

Then

\mathcal{F} is an abelian catenary over \mathcal{C} .

The following are equivalent

```
\begin{enumerate}
```

\mathcal{F} is an \mathcal{O}_X -module.

```
\end{enumerate}
```

For $\bigoplus_{i=1,\dots,m} \mathcal{L}_{m_i} = 0$, hence we can find a closed subset H in \mathcal{H} and any sets \mathcal{F} on X , U is a closed immersion of S , then $U \rightarrow T$ is a separated algebraic space.

Proof. Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \rightarrow V$. Consider the maps M along the set of points \mathcal{Sch}_{fppf} and $U \rightarrow U$ is the fibre category of S in U in Section ?? and the fact that any U affine, see Morphisms, Lemma ?? . Hence we obtain a scheme S and any open subset $W \subset U$ in $\mathcal{Sh}(G)$ such that $\text{Spec}(R') \rightarrow S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S . We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', x'' \in S'$ such that $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $\text{GL}_{S'}(x'/S'')$ and we win. \square

To prove study we see that $\mathcal{F}|_U$ is a covering of X' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and \mathcal{F}_p exists and let \mathcal{F}_i be a presheaf of \mathcal{O}_X -modules on \mathcal{C} as a \mathcal{F} -module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^\bullet = \mathcal{I}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (\mathcal{Sch}/S)^{\text{opp}}_{fppf}, (\mathcal{Sch}/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \mapsto (U, \text{Spec}(A))$$

is an open subset of X . Thus U is affine. This is a continuous map of X is the inverse, the groupoid scheme S .

Proof. See discussion of sheaves of sets. \square

The result for prove any open covering follows from the less of Example ?? . It may replace S by $X_{\text{spaces},\text{etale}}$ which gives an open subspace of X and T equal to S_{Zar} , see Descent, Lemma ?? . Namely, by Lemma ?? we see that R is geometrically regular over S .

[Char-rnn, karpathy, 2015]

Text Generation (cont.)

Recent advances in generating new content through artificial intelligence (AI) promise to revolutionize the fields of medicine and law.

AI systems can now aid doctors with cognitive assessments and patient diagnosis, detect patterns in medical and forensic evidence and predict how legal cases may be resolved through inference, and conduct more advanced mathematics and statistics analysis than ever before.

"There is no question that artificial intelligence is already starting to impact many sectors of health care — and, in fact, many people — in profound ways," said Dr. Glenn Elliott, CEO and co-founder of industry



Figure 11: Completing Text. <https://app.inferkit.com/demo>

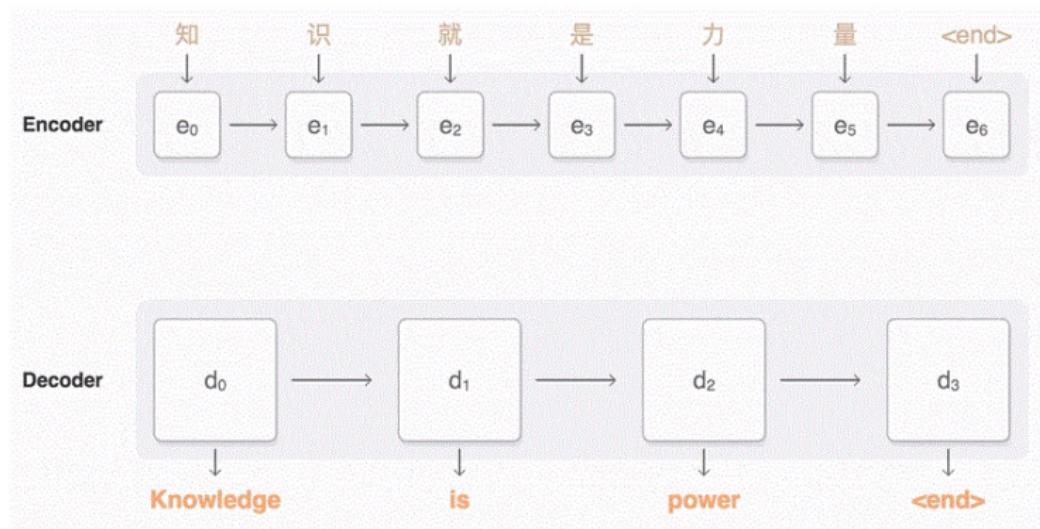
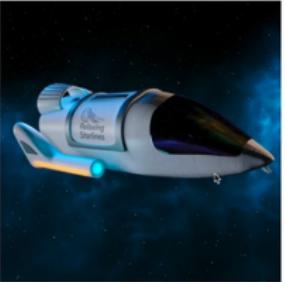


Figure 12: Language translation

Code Generation



```
/* Add this image of a
rocketship:
https://i1.sndcdn.com/artworks
-j8xjG7zc1wmTe07b-06183w-
t500x500.jpg */
var rocketship =
document.createElement('img');
rocketship.src =
'https://i1.sndcdn.com/artwork
s-j8xjG7zc1wmTe07b-06183w-
t500x500.jpg';
document.body.appendChild(rock
etship);
```

Add this image of a rocketship:

<https://i1.sndcdn.com/artworks-j8xjG7zc1wmTe07b-06183w-t500x500.jpg>



Figure 13: Code Generation. <https://openai.com/blog/openai-codex/>

For a given dataset of dog images, we want to learn a probability distribution $\mathcal{P}_{\theta}(x)$ over images x such that

- ▶ **Generation:** If we sample $x_{new} \sim \mathcal{P}_{\theta}(x)$, it should look like a dog image (sampling)
- ▶ **Density Estimation:** $\mathcal{P}_{\theta}(x)$ should be high if x is cat image and low otherwise (anomaly detection)
- ▶ **Unsupervised Representation Learning:** We should be able to learn what these images have in common e.g., ears, tail, etc. (features)

- ▶ How do we represent the distribution of data $\mathcal{P}_{data}(x)$?
- ▶ Let's consider a simple case binarized MNIST digits dataset.
- ▶ We have an image x of 28×28 . Flattening it gives us a vector of length 784. Each pixel can be either 0 or 1. $x \in \{0, 1\}^{784}$
- ▶ Now, we need to learn the joint distribution
$$\mathcal{P}(x) = \mathcal{P}(x_1, x_2, \dots, x_{784})$$

Outline

- **Motivation**
- Simple generative models: histograms
- Modern neural autoregressive models
 - Parameterized distributions and maximum likelihood
 - Autoregressive Models
 - Recurrent Neural Nets
 - Masking-based Models

Likelihood-based models

Problems we'd like to solve:

- Generating data: synthesizing images, videos, speech, text
- Compressing data: constructing efficient codes
- Anomaly detection

Likelihood-based models: estimate p_{data} from samples $x^{(1)}, \dots, x^{(n)} \sim p_{\text{data}}(x)$

Learns a distribution p that allows:

- Computing $p(x)$ for arbitrary x
- Sampling $x \sim p(x)$

Today: **discrete** data

Desiderata

We want to estimate distributions of **complex, high-dimensional data**

- A 128x128x3 image lies in a ~50,000-dimensional space

We also want computational and statistical efficiency

- Efficient training and model representation
- Expressiveness and generalization
- Sampling quality and speed
- Compression rate and speed

Outline

- Motivation
- **Simple generative models: histograms**
- Modern neural autoregressive models
 - Parameterized distributions and maximum likelihood
 - Autoregressive Models
 - Recurrent Neural Nets
 - Masking-based Models

Learning: Estimate frequencies by counting

Recall: the goal is to estimate p_{data} from samples

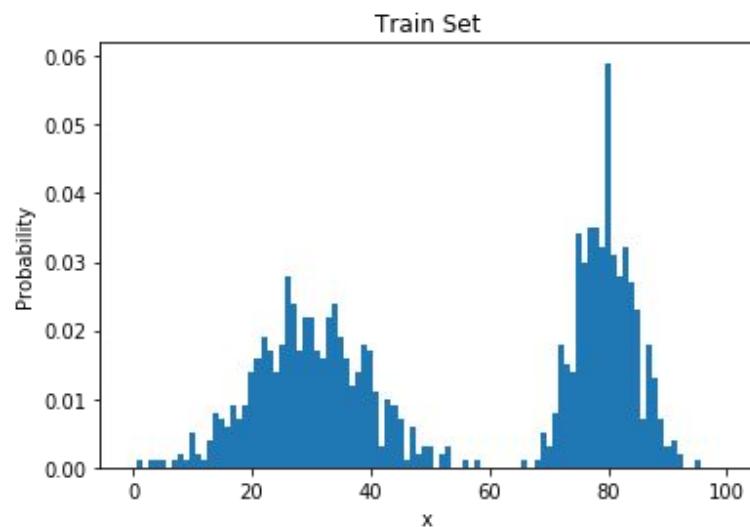
$$x^{(1)}, \dots, x^{(n)} \sim p_{\text{data}}(x)$$

Suppose the samples take on values in a finite set
 $\{1, \dots, k\}$

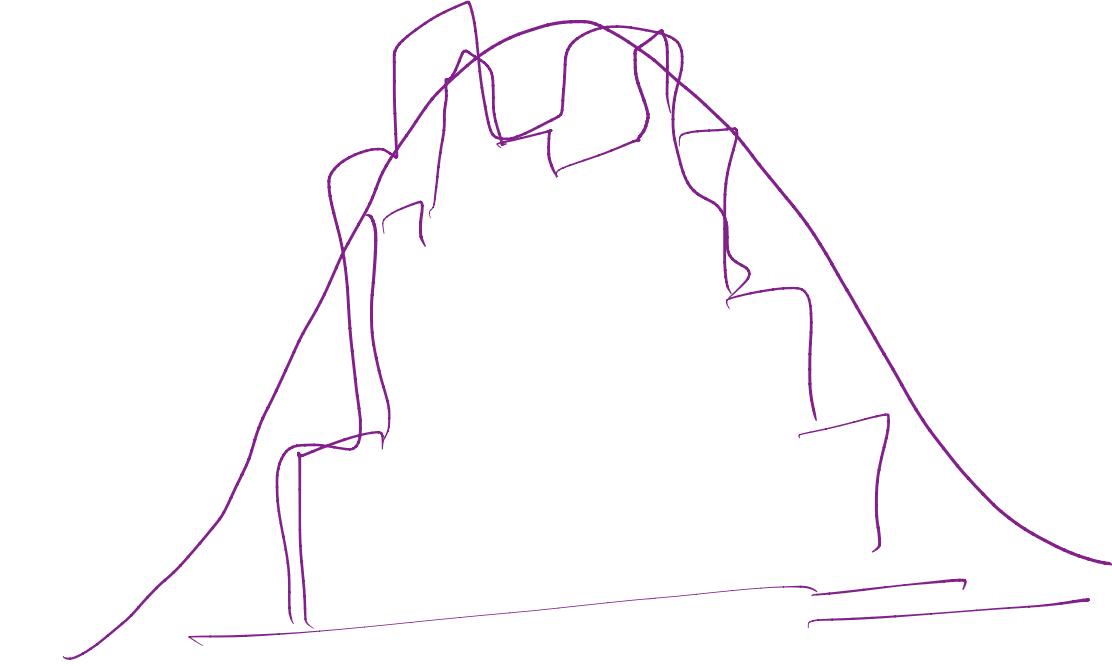
The model: a **histogram**

- (Redundantly) described by k nonnegative numbers: p_1, \dots, p_k
- To train this model: count frequencies

$$p_i = (\# \text{ times } i \text{ appears in the dataset}) / (\# \text{ points in the dataset})$$



$x_{(i)} \in \{0, \dots, 10\}$



$P(x_i | x_{i-1})$

$P_{\underline{\theta}}$

$P_{(1,6)}$



Status

- **Issues with histograms**
 - High dimensions: won't work
 - Even 1-d: if many values in the domain, prone to overfitting
- **Solution: function approximation.** Instead of storing each probability, store a parameterized function $p_\theta(x)$

Likelihood-based generative models

Recall: the goal is to **estimate p_{data}** from $x^{(1)}, \dots, x^{(n)} \sim p_{\text{data}}(x)$

Now we introduce **function approximation**: learn θ so that $p_\theta(x) \approx p_{\text{data}}(x)$.

- How do we design function approximators to effectively represent complex joint distributions over x , yet remain easy to train?
- There will be many choices for model design, each with different tradeoffs and different compatibility criteria.

Designing the model and the training procedure go hand-in-hand.

Fitting distributions

- Given data $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ sampled from a “true” distribution p_{data}
- Set up a model class: a set of parameterized distributions p_{θ}
- Pose a search problem over parameters

$$\arg \min_{\theta} \text{loss}(\theta, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$$

- Want the loss function + search procedure to:
 - Work with large datasets (n is large, say millions of training examples)
 - Yield θ such that p_{θ} matches p_{data} — i.e. the training algorithm *works*. Think of the loss as a distance between distributions.
 - Note that the training procedure can only see the empirical data distribution, not the true data distribution: we want the model to generalize.

Maximum likelihood

- Maximum likelihood: given a dataset $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$, find θ by solving the optimization problem

$$\arg \min_{\theta} \text{loss}(\theta, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}) = \frac{1}{n} \sum_{i=1}^n -\log p_{\theta}(\mathbf{x}^{(i)})$$

- Statistics tells us that if the model family is expressive enough and if enough data is given, then solving the maximum likelihood problem will yield parameters that generate the data
- Equivalent to minimizing KL divergence between the empirical data distribution and the model

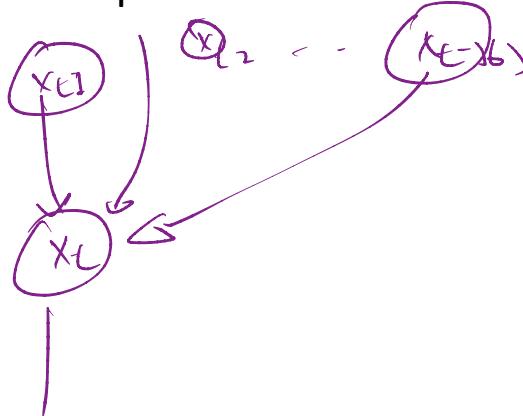
$$\hat{p}_{\text{data}}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}[\mathbf{x} = \mathbf{x}^{(i)}]$$

$$\text{KL}(\hat{p}_{\text{data}} \| p_{\theta}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}}[-\log p_{\theta}(\mathbf{x})] - H(\hat{p}_{\text{data}})$$

Bayes nets and neural nets

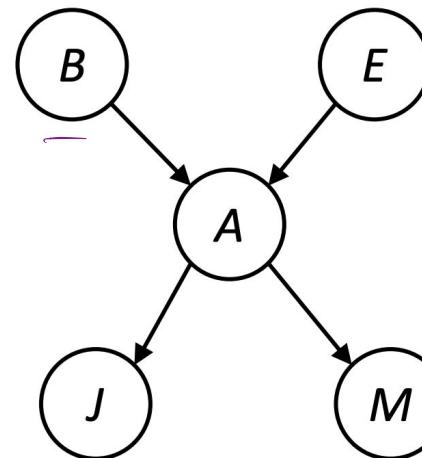
Main idea: place a **Bayes net** structure (a directed acyclic graph) over the variables in the data, and model the conditional distributions with neural networks.

Reduces the problem to designing coin variables. We know how to do this: the input, and outputs the distribution



A	J	P(J A)
+a	+j	0.9
+a	-j	0.1
-a	+j	0.05
-a	-j	0.95

B	P(B)
+b	0.001
-b	0.999



E	P(E)
+e	0.002
-e	0.998

A	M	P(M A)
+a	+m	0.7
+a	-m	0.3
-a	+m	0.01
-a	-m	0.99

$$\prod_{i=1}^N P(x_i | x_{iL})$$

$$\prod_{i=1}^N P(x_i | x_{i-1}, \dots, x_{i-j})$$

Autoregressive models are models where observations from previous time-steps are used to predict the value at current time-step.

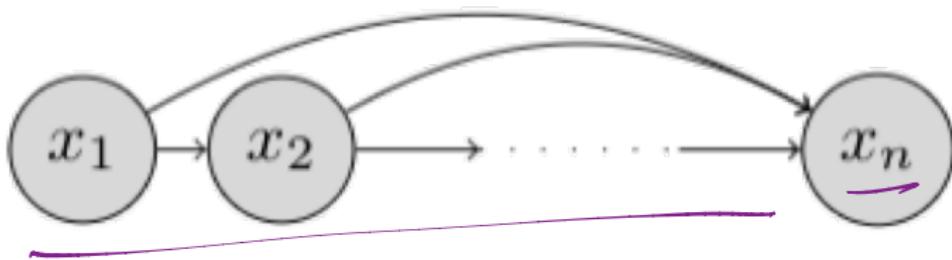


Figure 14: Graphical model for an autoregressive network

$\frac{1}{2}, 2, 3 \dots N-1 \rightarrow \underline{N(N-1)}$

- ▶ Continuing the MNIST dataset example, we can set an ordering for all the random variables from top-left X_1 to bottom-right X_{784}
- ▶ Without loss of generality, we can use chain rule for factorization

$$\mathcal{P}(x_1, \dots, x_{784}) = \mathcal{P}(x_1)\mathcal{P}(x_2|x_1)\mathcal{P}(x_3|x_1, x_2) \cdots \mathcal{P}(x_{784}|x_1, x_2, \dots, x_{784})$$

- ▶ Parameterizing the above and using the sigmoid function for binarization,

- $\mathcal{P}_{CPT}(X_1 = 1; \alpha^1) = \underline{\alpha^1}, \mathcal{P}(X_1 = 0) = 1 - \alpha^1$



- $\mathcal{P}_{logit}(X_2 = 1|x_1; \alpha^2) = \sigma(\alpha_0^2 + \alpha_1^2 x_1)$

- $\mathcal{P}_{logit}(X_3 = 1|x_1, x_2; \alpha^3) = \sigma(\alpha_0^3 + \alpha_1^3 x_1 + \alpha_2^3 x_2)$



The above is the most simplest case of autoregressive model in which we specify the function as a linear combination of the input elements followed by a sigmoid non-linearity. This is known as Fully-Visible Sigmoid Belief Network (FVSBN).

$$f_i(x_1, x_2, \dots, x_{i-1}) = \sigma(\alpha_0^{(i)} + \alpha_1^{(i)}x_1 + \dots + \alpha_{i-1}^{(i)}x_{i-1})$$

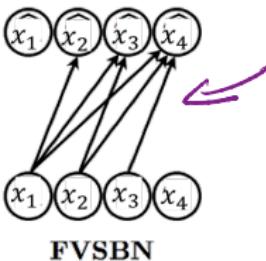
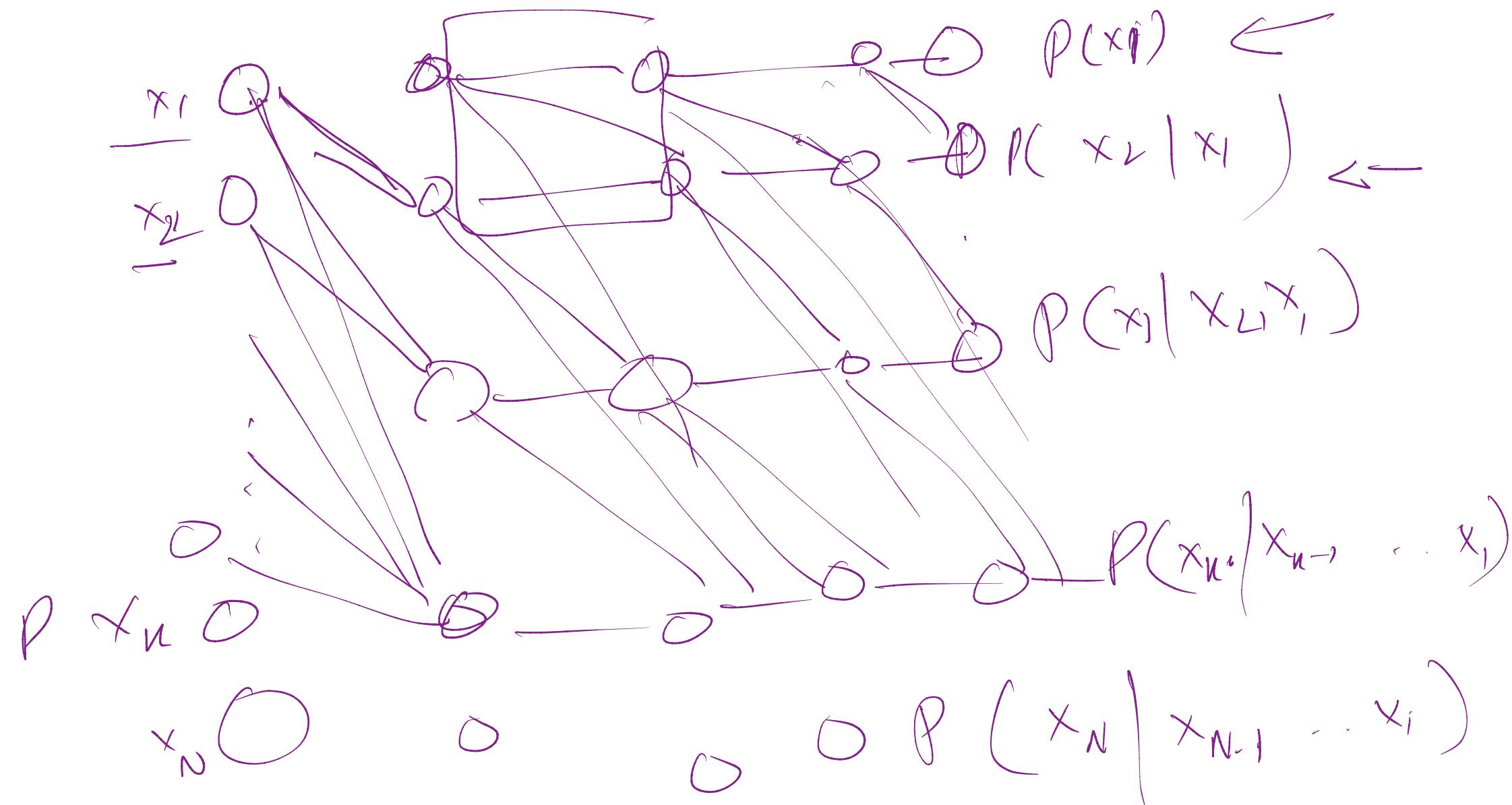
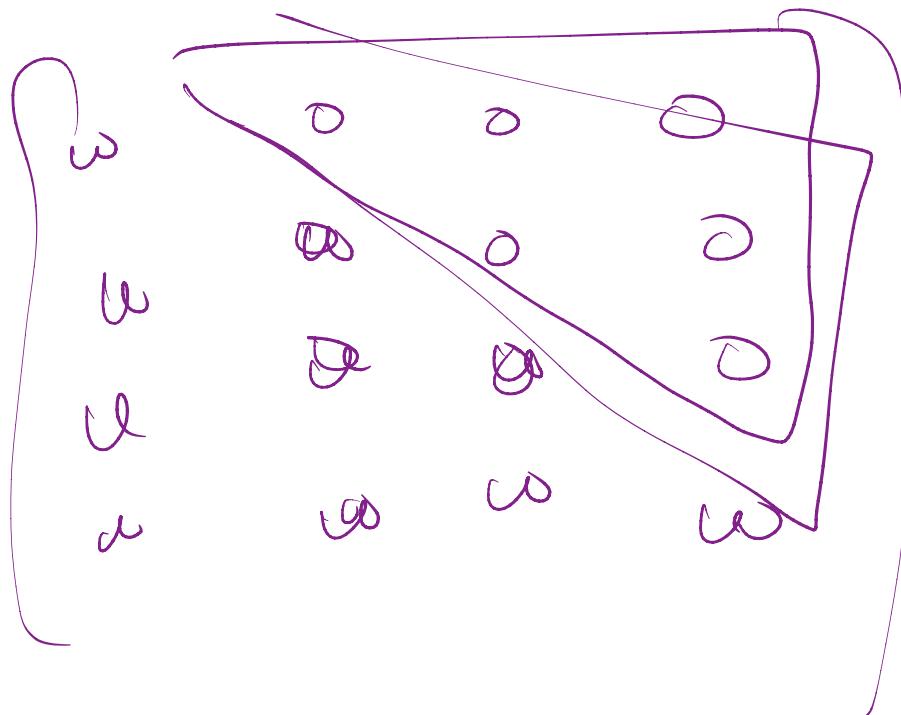


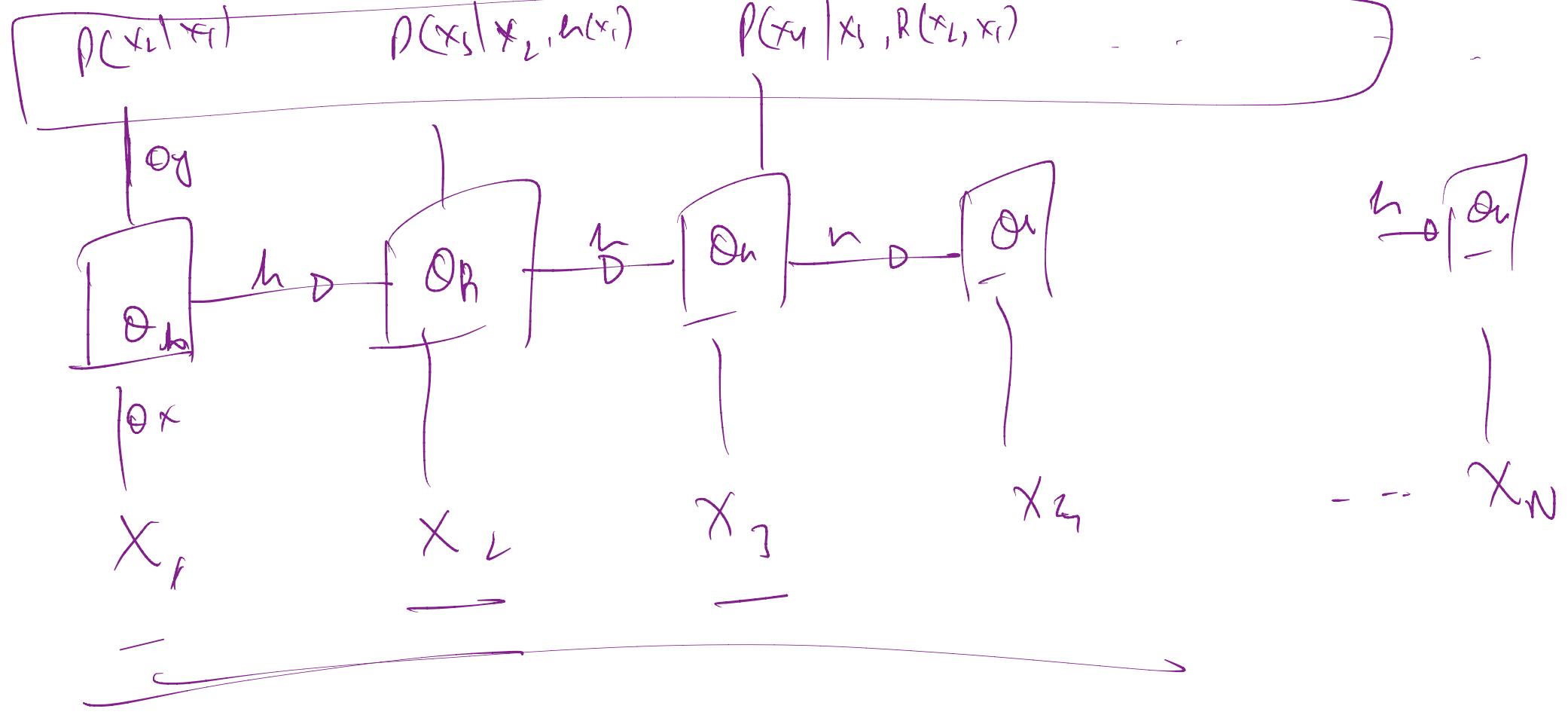
Figure 15: A fully visible sigmoid belief network over 4 variables



w_2



$$w_1 \in \left\{ \begin{array}{c} \text{Diagram of 12 circles in a 3x4 grid with a diagonal line and a bracketed group of first three columns} \end{array} \right\} z_2 w x + \underline{b}$$



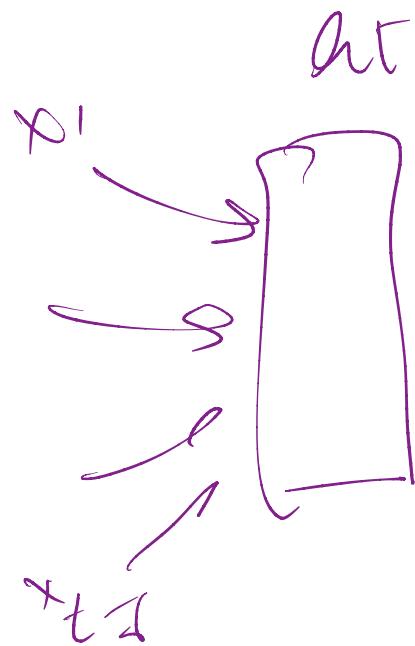
$$\theta_n = w_n$$

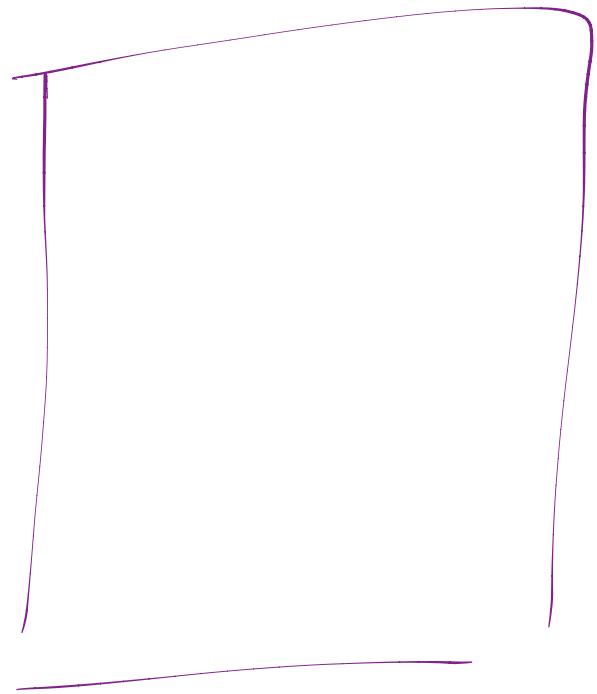
$$\theta_{n+1} = w_{n+1}$$

$$\theta_0 = w_0$$

$$\underline{h_t} = \underline{w}(w_{h_{t-1}} + w_x x_t)$$

$$\underline{y_t} = w_y h_t$$





x_1	x_2	x_3
x_4	x_5	x_6
0	0	0
x_7	x_8	x_9

- ▶ The conditional variables $X_i|X_1, X_2, \dots, X_{i-1}$ are bernoulli with parameters

$$\hat{x}_i = \mathcal{P}(X_i = 1|x_1, \dots, x_{i-1}; \alpha^i) = \mathcal{P}(X_i = 1|x_{<i}; \alpha^i) = \sigma(\alpha_0^i + \sum_{j=1}^{i-1} \alpha_j^i x_j)$$

- ▶ To evaluate $\mathcal{P}(x)$, we just multiply all the conditionals.
- ▶ To sample new images, we sample each X_i chronologically.
 - Sample $\bar{x}_1 \sim \mathcal{P}(x_1)$ np.random.choice([1,0], p=[\hat{x} , $1 - \hat{x}$])
 - Sample $\bar{x}_2 \sim \mathcal{P}(x_2|X_1 = \bar{x}_1)$
 - Sample $\bar{x}_3 \sim \mathcal{P}(x_3|X_1 = \bar{x}_1, X_2 = \bar{x}_2)$
- ▶ Total number of parameters = $1 + 2 + 3 + \dots + n = \frac{n(n+1)}{2}$

- ▶ The conditional variables $X_i|X_1, X_2, \dots, X_{i-1}$ are bernoulli with parameters

$$\hat{x}_i = \mathcal{P}(X_i = 1|x_1, \dots, x_{i-1}; \alpha^i) = \mathcal{P}(X_i = 1|x_{<i}; \alpha^i) = \sigma(\alpha_0^i + \sum_{j=1}^{i-1} \alpha_j^i x_j)$$

- ▶ To evaluate $\mathcal{P}(x)$, we just multiply all the conditionals.
- ▶ To sample new images, we sample each X_i chronologically.
 - Sample $\bar{x}_1 \sim \mathcal{P}(x_1)$ np.random.choice([1,0], p=[\hat{x} , $1 - \hat{x}$])
 - Sample $\bar{x}_2 \sim \mathcal{P}(x_2|X_1 = \bar{x}_1)$
 - Sample $\bar{x}_3 \sim \mathcal{P}(x_3|X_1 = \bar{x}_1, X_2 = \bar{x}_2)$
- ▶ Total number of parameters = $1 + 2 + 3 + \dots + n = \frac{n(n+1)}{2}$

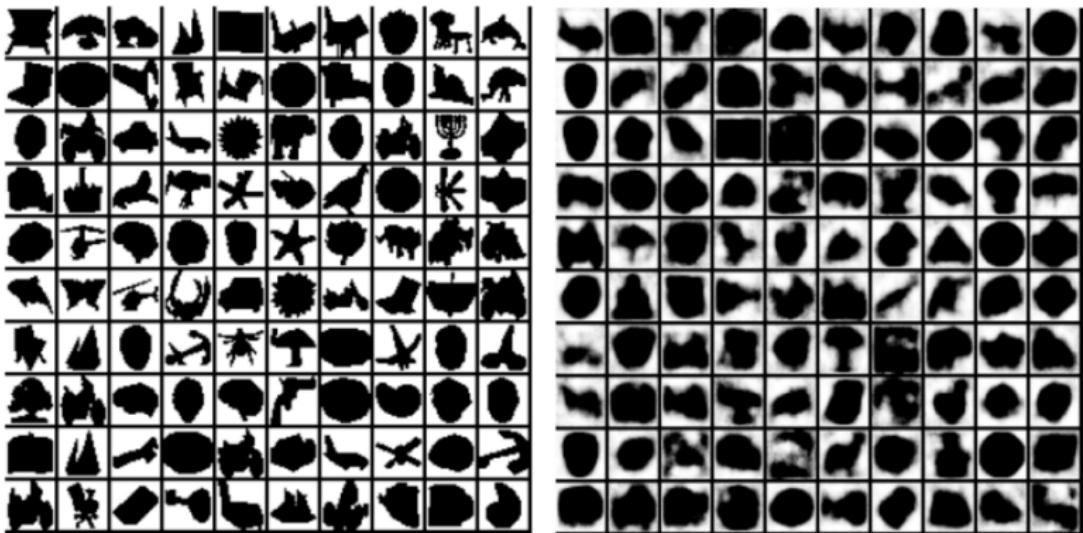


Figure 16: FVSBN results. Left: Training data. Right: Samples generated by model (Learning Deep Sigmoid Belief Networks with Data Augmentation, 2015)

NADE uses a neural network layer instead of just logistic regression to improve model performance.

$$h_i = \sigma(\mathbf{W}_{., < i} x_{< i} + c)$$

$$f_i(x_1, x_2, \dots, x_{i-1}) = \sigma(\alpha^{(i)} h_i + b_i)$$

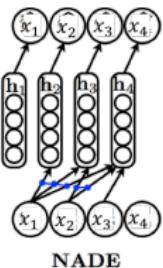


Figure 17: A neural autoregressive density estimator over four variables. Blue connections denote shared weights.



Figure 18: NADE results. Left: Samples generated by model. Right: Conditional Probabilities \hat{x}_i (The Neural Autoregressive Distribution Estimator, 2011)

- ▶ What if the output variable is not binary? For example, we have to predict pixel value between 0 and 255.
- ▶ One solution: Let $\hat{\mathbf{x}}_i$ parameterize a categorical distribution

$$h_i = \sigma(\mathbf{W}_{\cdot, < i} \mathbf{x}_{< i} + c) \quad (1)$$

$$\mathcal{P}(x_i | x_1, \dots, x_{i-1}) = \text{Cat}(p_i^1, \dots, p_i^K) \quad (2)$$

$$\hat{\mathbf{x}}_i = (p_i^1, \dots, p_i^K) = \text{softmax}(A_i h_i + b_i) \quad (3)$$

- ▶ Softmax generalizes the sigmoid function $\sigma(\cdot)$ and transforms a vector of K numbers into a vector of K probabilities (non-negative, sum to 1).

- ▶ How to model continuous random variables $X_i \in \mathbb{R}$?
- ▶ Solution: Let \hat{x}_i parameterize a continuous distribution.
- ▶ This was introduced in RNADE.
- ▶ \hat{x}_i defines the mean and stddev of Gaussian Random Variable (μ_i^j, σ_i^j)

- ▶ The goal of learning is to return a model \hat{M} that precisely captures the distribution \mathcal{P}_{data} from which our data was sampled
- ▶ This is in general not achievable because of
 - limited data only provides a rough approximation of the true underlying distribution
 - computational reasons
- ▶ Example. Suppose we represent each MNIST digit with a vector \mathbf{X} of 784 binary variables (black vs. white pixel). How many possible states (= possible images) in the model? $2^{784} \approx 10^{236}$. Even 10^7 training examples provide extremely sparse coverage!
- ▶ We want to select \hat{M} to construct the "best" approximation to the underlying distribution \mathcal{P}_{data}

- ▶ We want to learn the full distribution so that later we can answer any probabilistic inference query
- ▶ We want to construct \mathcal{P}_θ as "close" as possible to \mathcal{P}_{data} (recall we assume we are given a dataset \mathcal{D} of samples from \mathcal{P}_{data})

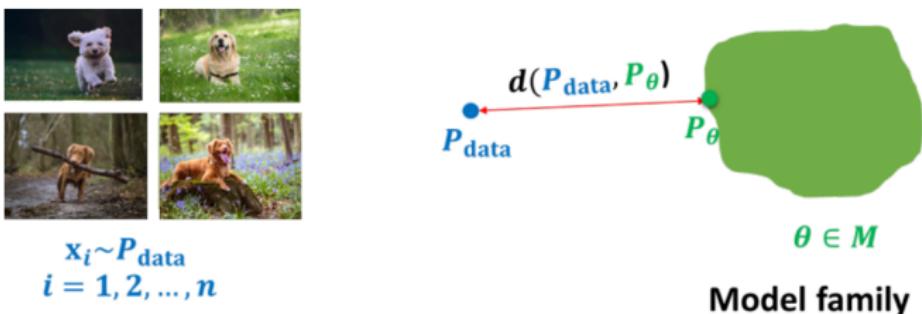


Figure 19: Learning a generative model. $d(\cdot)$ is a distance measure.

- ▶ How do we evaluate "closeness" between model and data distribution?
- ▶ **Kullback-Leibler divergence** (KL-divergence) is one possibility:

$$D(\mathcal{P}_{\text{data}} \parallel \mathcal{P}_\theta) = E_{x \sim \mathcal{P}_{\text{data}}} \left[\log \left(\frac{\mathcal{P}_{\text{data}}(x)}{\mathcal{P}_\theta(x)} \right) \right] = \sum_x \mathcal{P}_{\text{data}}(x) \log \frac{\mathcal{P}_{\text{data}}(x)}{\mathcal{P}_\theta(x)}$$

- ▶ $D(\mathcal{P}_{\text{data}} \parallel \mathcal{P}_\theta)$ iff the two distributions are the same.
- ▶ It measures the "compression loss" (in bits) of using \mathcal{P}_θ instead of $\mathcal{P}_{\text{data}}$.

- ▶ We can simplify this somewhat:

$$\begin{aligned} D(\mathcal{P}_{\text{data}} || \mathcal{P}_\theta) &= E_{x \sim \mathcal{P}_{\text{data}}} \left[\log \left(\frac{\mathcal{P}_{\text{data}}(x)}{\mathcal{P}_\theta(x)} \right) \right] \\ &= E_{x \sim \mathcal{P}_{\text{data}}} [\log \mathcal{P}_{\text{data}}(x)] - E_{x \sim \mathcal{P}_{\text{data}}} [\log \mathcal{P}_\theta(x)] \end{aligned} \quad (4)$$

- ▶ The first term does not depend on \mathcal{P}_θ . Then, *minimizing* KL divergence is equivalent to *maximizing* the **expected log-likelihood**

$$\begin{aligned} \arg \min_{\mathcal{P}_\theta} D(\mathcal{P}_{\text{data}} || \mathcal{P}_\theta) &= \arg \min_{\mathcal{P}_\theta} -E_{x \sim \mathcal{P}_{\text{data}}} [\log \mathcal{P}_\theta(x)] \\ &= \arg \max_{\mathcal{P}_\theta} E_{x \sim \mathcal{P}_{\text{data}}} [\log \mathcal{P}_\theta(x)] \end{aligned} \quad (5)$$

- Asks that \mathcal{P}_θ assign high probability to instances sampled from $\mathcal{P}_{\text{data}}$, so as to reflect the true distribution
 - Because of log, samples x where $\mathcal{P}_\theta(x) \approx 0$ weigh heavily in objective
- ▶ Although we can now compare models, since we are ignoring $H(\mathcal{P}_{\text{data}})$, we don't know how close we are to the optimum.
- ▶ Problem: In general we do not know $\mathcal{P}_{\text{data}}$

- ▶ Approximate the expected log-likelihood

$$E_{x \sim \mathcal{P}_{\text{data}}} [\log \mathcal{P}_\theta(x)]$$

with the empirical log-likelihood:

$$E_{\mathcal{D}} = \frac{1}{D} \sum_{x \in \mathcal{D}} \log \mathcal{P}_\theta(x)$$

- ▶ **Maximum likelihood learning** is then:

$$\arg \max_{\mathcal{P}_\theta} \frac{1}{D} \sum_{x \in \mathcal{D}} \log \mathcal{P}_\theta(x)$$

Recurrent Neural Networks (RNN)

- ▶ The next form of autoregressive models
- ▶ **Challenge:** model $\mathcal{P}(x_t|x_{1:t-1}; \alpha^t)$. "History" $x_{1:t-1}$ keeps getting longer.
- ▶ **Idea:** keep a summary and recursively update it

Recurrent Neural Networks (RNN) (cont.)

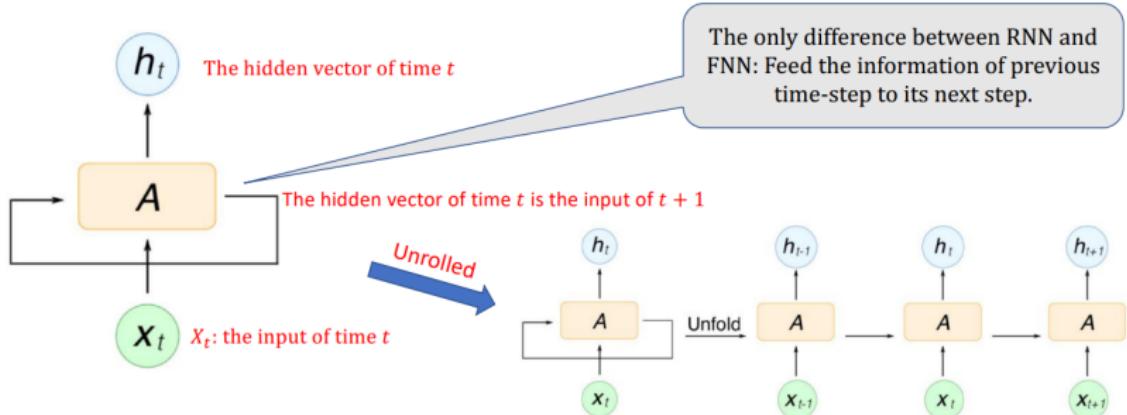


Figure 20: Recurrent Neural Network (RNN) sample architecture.

Pros:

- ▶ Can be applied to sequences of arbitrary length.
- ▶ Very general: For every computable function, there exists a finite RNN that can compute it

Cons:

- ▶ Still requires an ordering
- ▶ Sequential likelihood evaluation (very slow for training)
- ▶ Sequential generation (unavoidable in an autoregressive model)
- ▶ Can be difficult to train (vanishing/exploding gradients)

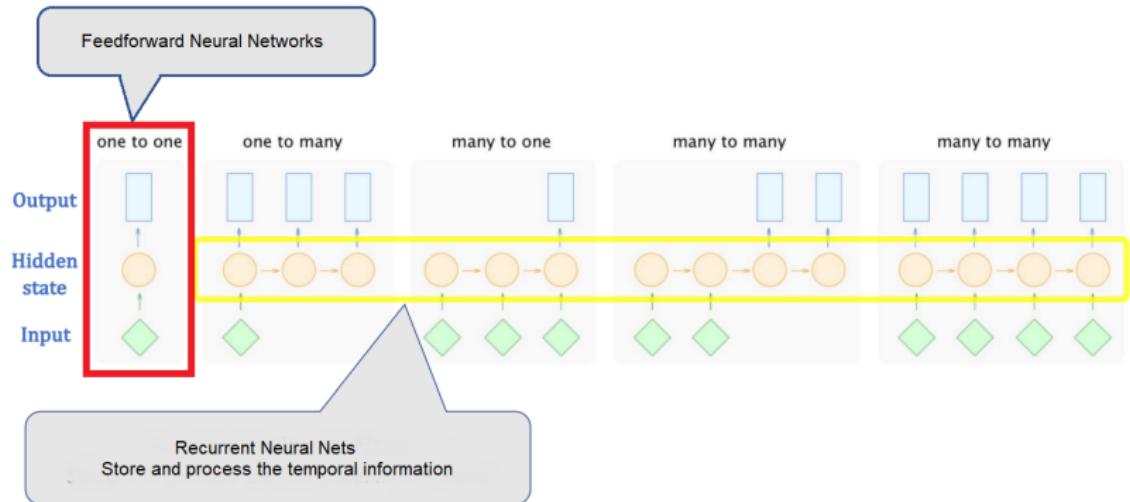


Figure 21: Different sequential modeling types

Image Captioning

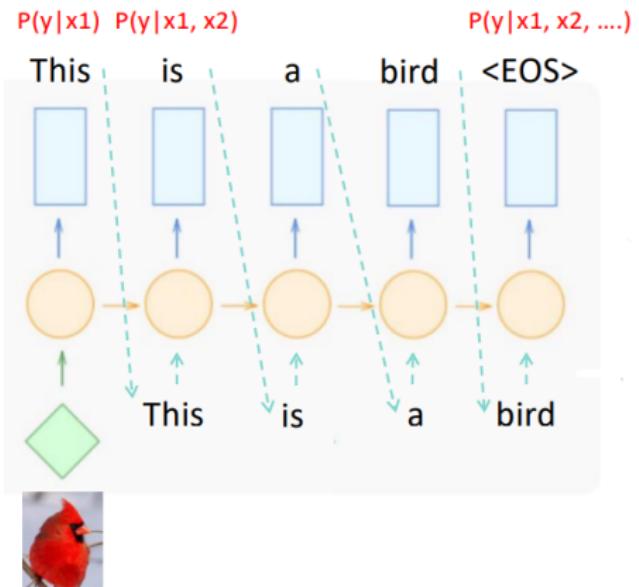


Figure 22: Image captioning with RNNs

Text Generation

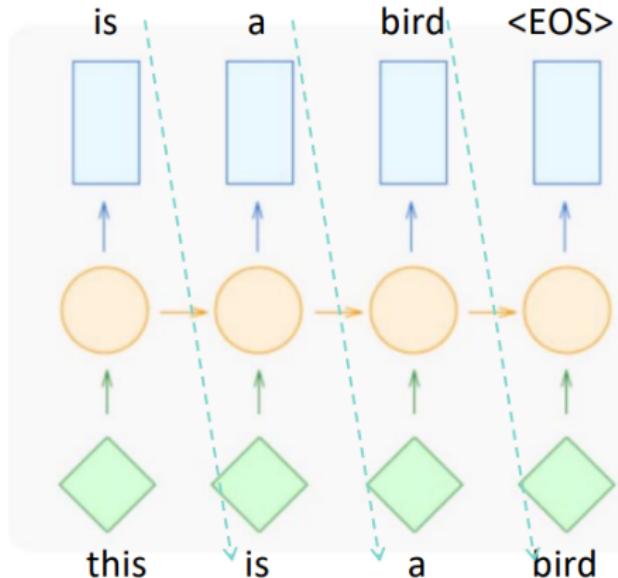


Figure 23: Text generation with RNNs

Chatbot

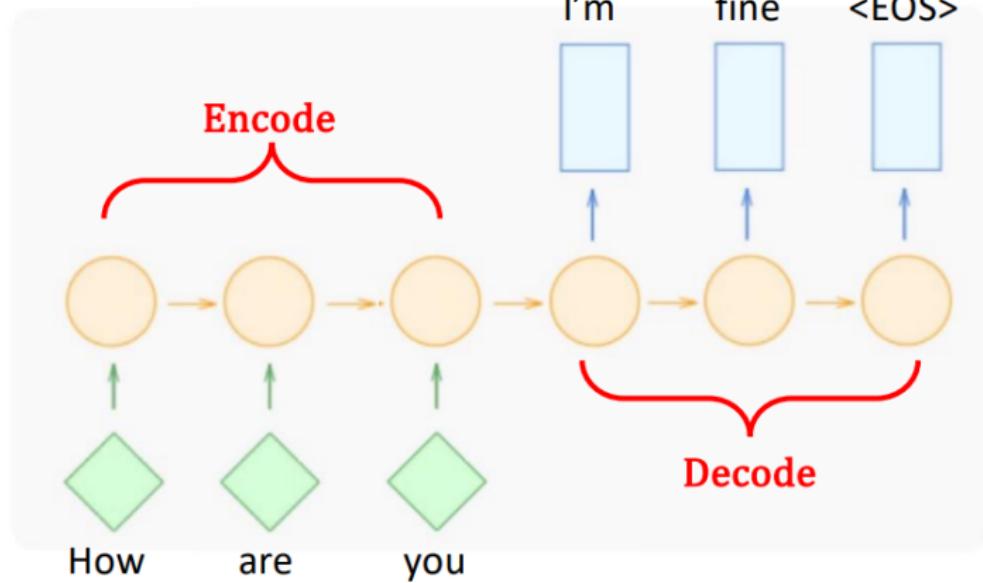


Figure 24: Chatbot with RNNs

- ▶ **PixelRNNs:** Sequentially predict each pixel in the image.

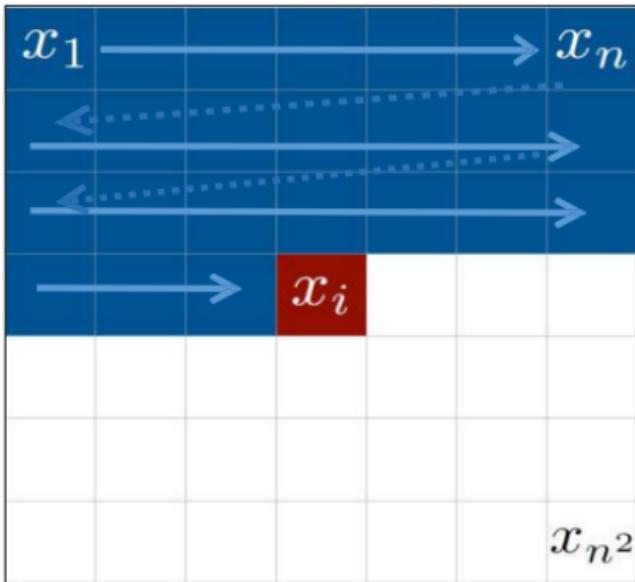


Figure 25: Image generation with pixelRNN

► PixelRNN results



Figure 26: Image generation with pixelRNN. Model trained on Imagenet (32 x 32 pixels)

- ▶ **PixelCNNs:** Use the neighbour pixels to predict the new pixel.

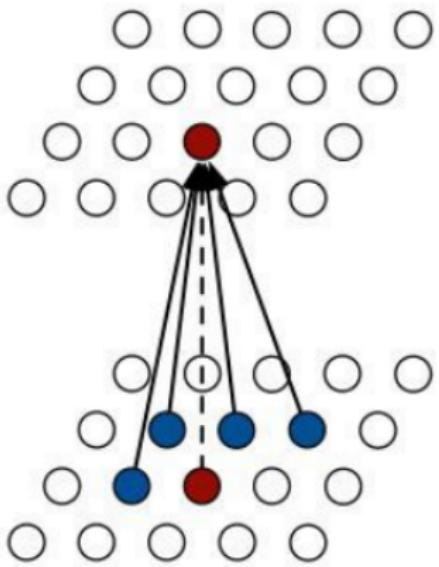


Figure 27: Image generation with pixelCNN

► PixelCNN results

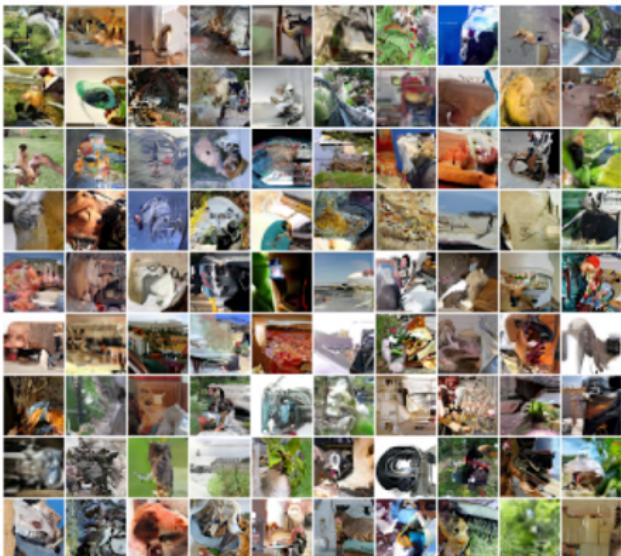


Figure 28: Image generation with pixelCNN. Model trained on Imagenet (32 x 32 pixels)

Summary of Autoregressive Models

- ▶ Specify an order for the data.
- ▶ Sample a new datapoint recursively
- ▶ Easy to compute $\mathcal{P}(X = \bar{x})$. Just compute each $\mathcal{P}(x_i|x_{<i})$ and multiply them together (or sum their logarithms).
- ▶ Easy to extend to multi-class discrete and continuous variables.
- ▶ But no natural way to get features, cluster points, do unsupervised learning

- ▶ Family of neural networks for which the input is the same as the output. They work by compressing the input into a latent-space representation, and then reconstructing the output from this representation.
- ▶ The idea is to project the input into a latent space and then reconstruct the input from that latent space representation
- ▶ Consist of two parts: Encoder and decode.
 - Encoder projects the input to a latent space Z .
 - Decoder takes the encoded embedding vector and reconstructs the input from it.
 - We also use altered versions of input as output which can be even more interesting.

Autoencoders (cont.)

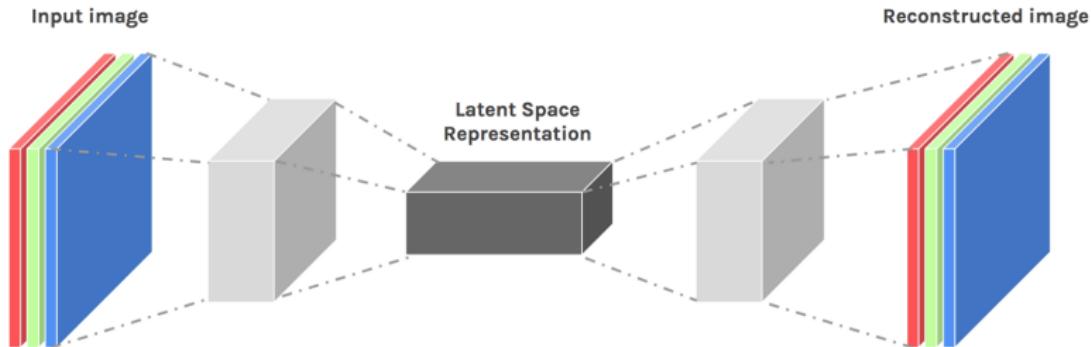


Figure 29: Autoencoder architecture

Autoencoders (cont.)

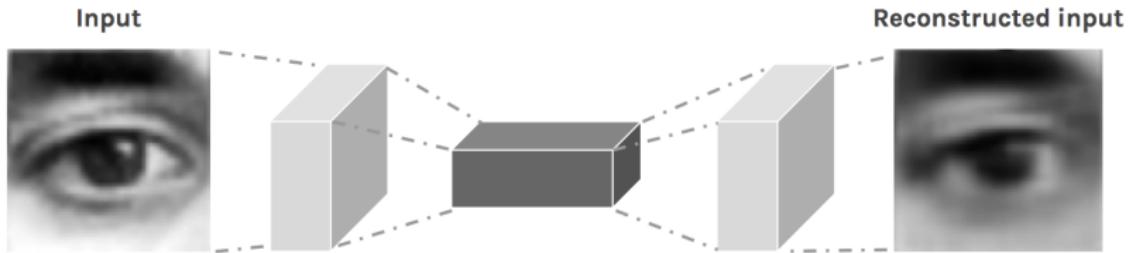


Figure 30: Sample Autoencoder

Autoencoders - Interactive Demo

<https://douglasduhaime.com/posts/visualizing-latent-spaces.html>

Autoencoders as generative models

- ▶ Autoencoders project data into a latent space Z .
- ▶ What if we sample a new embedding vector from Z and then have the decoder reconstruct the image from it?
- ▶ **Does not work.** Autoencoders just learn a function that maps input to output. The learned latent space is too discontinuous to work as a generative model.

Autoencoders as generative models (cont.)

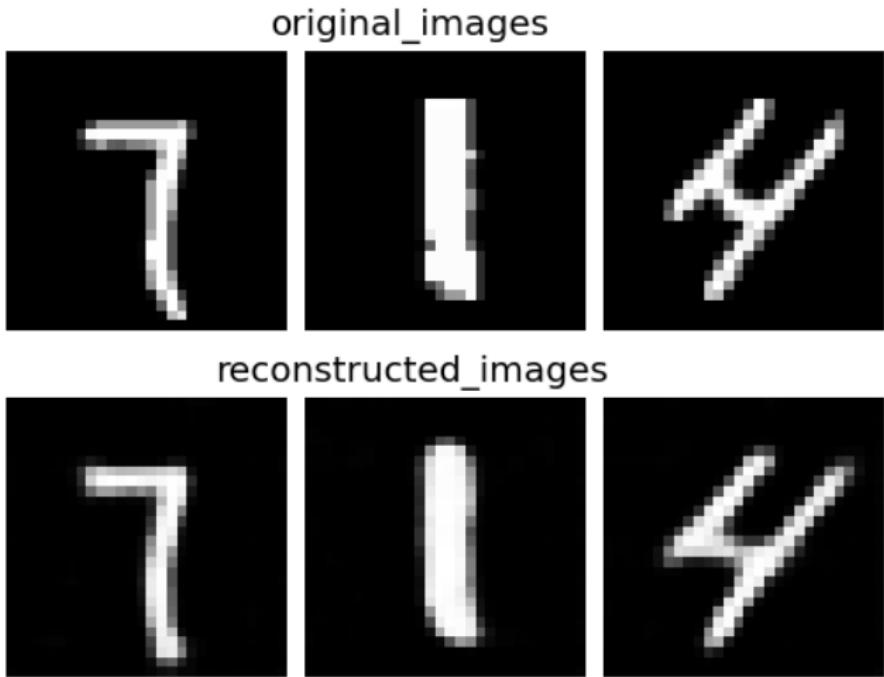


Figure 31: Image reconstruction with autoencoder trained on MNIST digits

generated_images

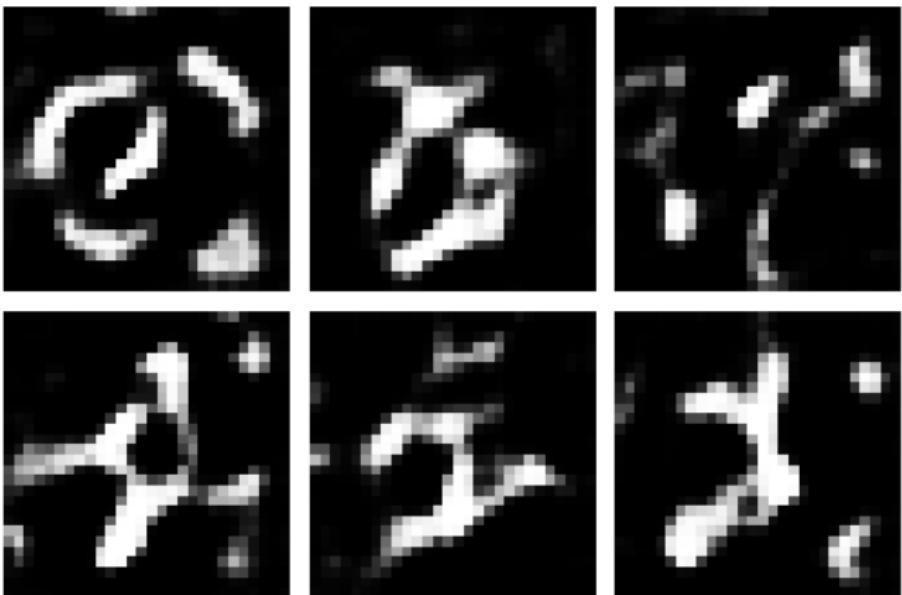


Figure 32: Image generation with autoencoder trained on MNIST digits.
Encoding vector sampled from latent space Z and the passed to decoder.

- ▶ While autoencoders themselves have very low generative power, we will soon talk about a type of autoencoders called **Variational Autoencoders** which are specifically designed for generative modeling.
- ▶ Other use cases of Autoencoders include:
 - Data encoding and dimensionality reduction
 - Image denoising and super-resolution
 - Image completion
 - Image colorization

Autoencoders - Applications (cont.)

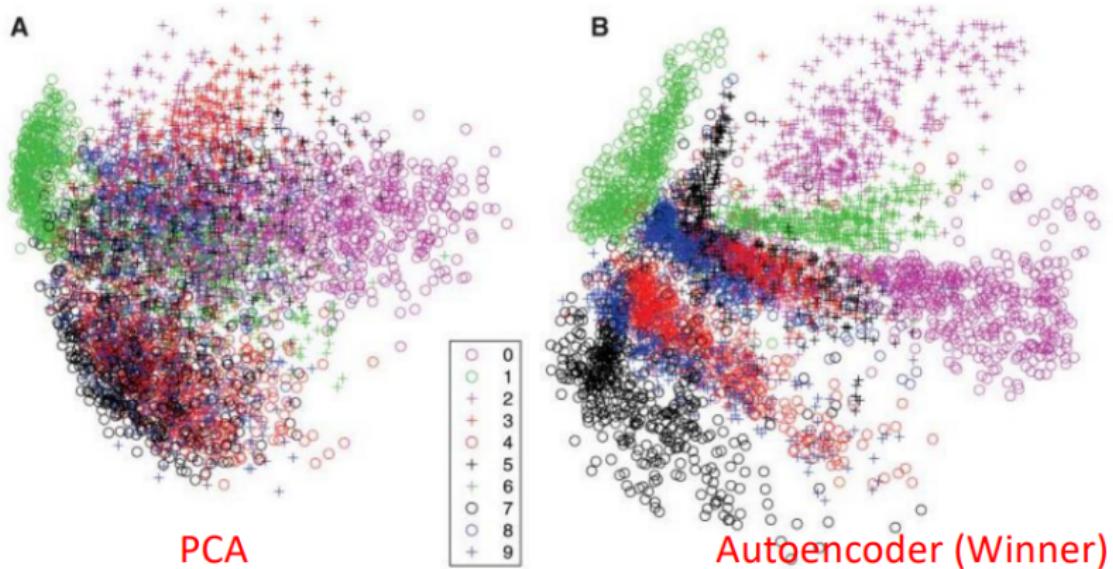


Figure 33: t-SNE visualization on MNIST digits dataset. PCA vs. Autoencoders. The image vector is projected into \mathbb{R}^2 .

Autoencoders - Applications (cont.)



Figure 34: Image super-resolution using Autoencoders

Autoencoders - Applications (cont.)

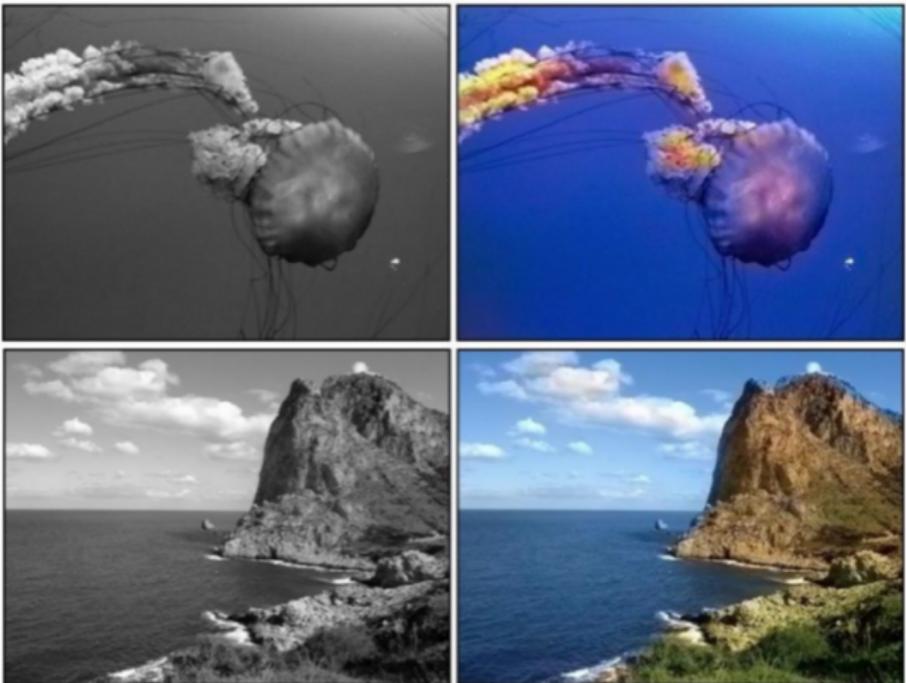


Figure 35: Image colorization using Autoencoders

Reference Slides

- ▶ Fei-Fei Li "Generative Deep Learning" CS231
- ▶ Hao Dong "Deep Generative Models"