

Title: Machine Learning (2021/2022)

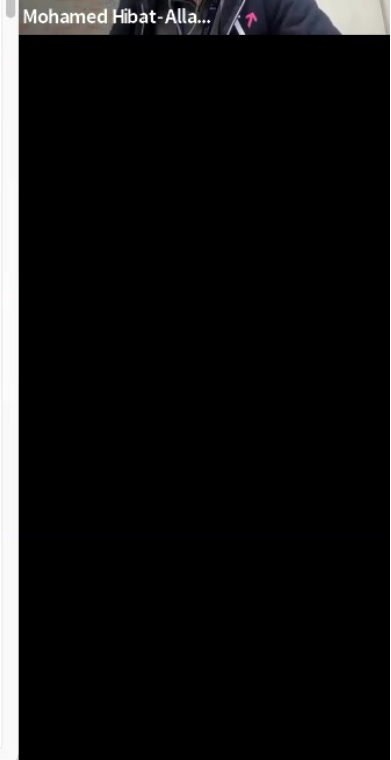
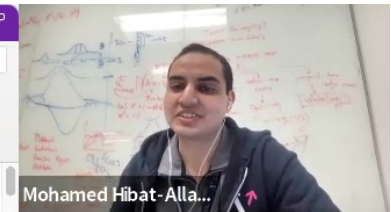
Speakers: Lauren Hayward

Collection: Machine Learning (2021/2022)

Date: April 26, 2022 - 11:30 AM

URL: <https://pirsa.org/22040074>

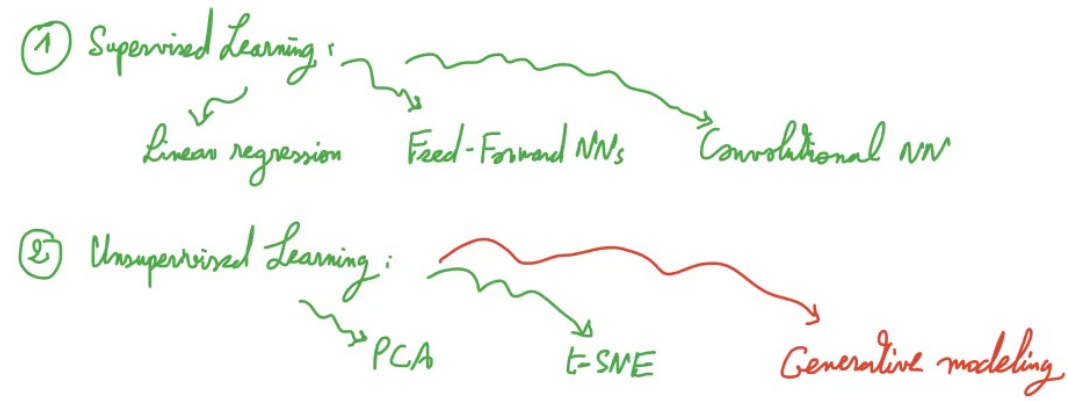
Abstract: This course is designed to introduce modern machine learning techniques for studying classical and quantum many-body problems encountered in condensed matter, quantum information, and related fields of physics. Lectures will focus on introducing machine learning algorithms and discussing how they can be applied to solve problem in statistical physics. Tutorials and homework assignments will concentrate on developing programming skills to study the problems presented in lecture.



Introduction to generative modeling

Introduction to generative modeling

↳ Review:



↳ Outline:

OneNote

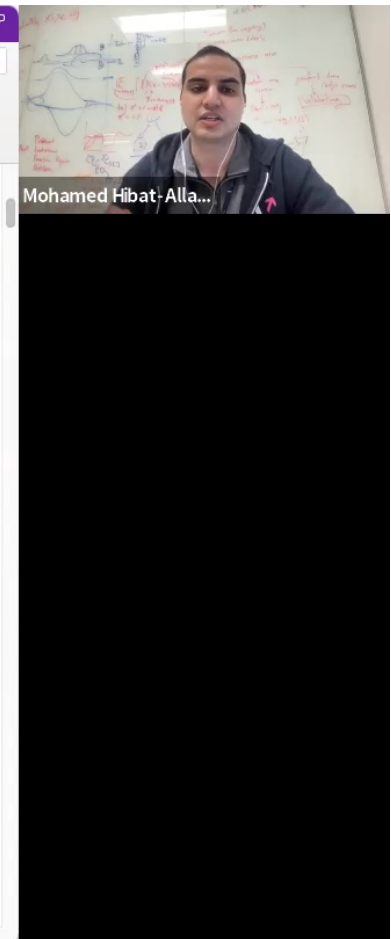
Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 2 mm 3 mm 5 mm 6.5 mm 8 mm

↪ Outline:

- 1 - Motivation for generative modeling:
- 2 - Maximum likelihood Principle & KL divergence:
- 3 - Classification of generative models
 - ↪ 3.1 - Exact likelihood models:
 - ↪ 3.2 - Approximate likelihood models:
 - ↪ 3.3 - Implicit likelihood models:

⊕ Motivation to generative modeling:



OneNote


Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 0.25 mm 0.35 mm 0.5 mm 0.7 mm 1 mm

→ 3.2 - Approximate likelihood models:

→ 3.3 - Implicit likelihood models:

① Motivation to generative modeling:



I am AI - AI Composed...

MH




Mohamed Hibat-Alla...

OneNote



Home Insert Draw View **Audio** Class Notebook

STANDING BY Status Level Record Stop Play Pause 00:32 / 02:48 Back 15 Seconds Forward 15 Seconds Add Bookmark

MH



Credit: this-person-does-not-exist.com



Mohamed Hibat-Alla...

OneNote

Home Insert Draw View **Audio** Class Notebook

STANDING BY Status

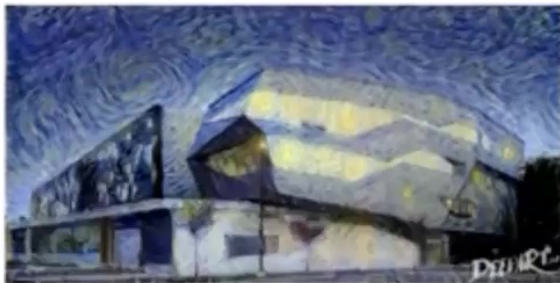
Level Record Stop

Play Pause 00:32 / 02:48

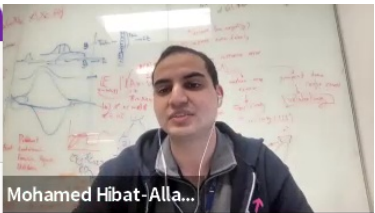
Back 15 Seconds Forward 15 Seconds

Add Bookmark

Share



Credit: Lauren Hayward, deepart.io



Mohamed Hibat-Alla...




MH

OneNote

Home Insert Draw View **Audio** Class Notebook


STANDING BY Status Level Record Stop Play Pause 00:32 / 02:48 Back 15 Seconds Forward 15 Seconds Add Bookmark

Credit: Lauren Hayward, deepart.io



Credit: Deep Dream Generator

MH

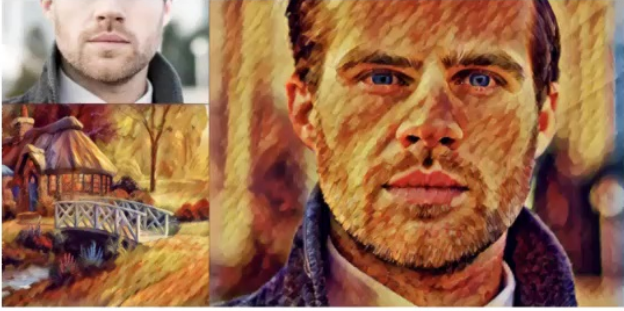


Mohamed Hibat-Alla...

OneNote

Home Insert Draw View **Audio** Class Notebook

STANDING BY Status Level Record Stop Play Pause 00:32 / 02:48 Back 15 Seconds Forward 15 Seconds Add Bookmark



Credit: Deep Dream Generator


Machine learning for many-body physics course

The machine learning course is a great way to learn about machine learning without having to learn a lot of math. It covers the basics of machine learning, and it's a great way to get started with machine learning.

The course has a lot of videos that you can watch on your own time. You can also download the course materials to work through the material at your own pace.

Credit: <https://app.inferkit.com/demo>

MH



Mohamed Hibat-Alla...

OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color

0.25 mm 0.35 mm 0.5 mm 0.7 mm 1 mm

can also download the course materials to work through the material at your own pace.

Credit: <https://app.inferkit.com/demo>

* The goal:

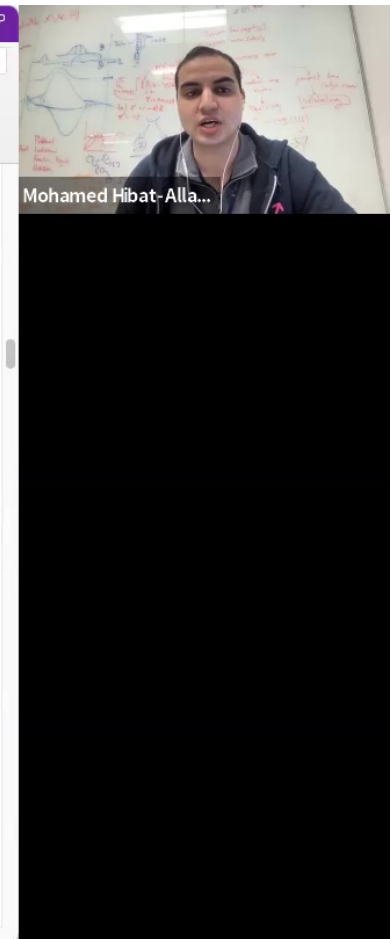
↳ We have input training data from some unknown distribution « P_{data} » and we want to have a model so that $P_{model} \approx P_{data}$.

D_{train} → Model → Training → $P_{model} \approx P_{data}$

$(D_{train} \sim P_{data})$
 ✓ ?

Sample new data points \vec{x}

MH



OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 0.25 mm 0.35 mm 0.5 mm 0.7 mm 1 mm

D_{train} → Model → Training → $P_{model} \approx P_{data}$
 ($D_{train} \approx P_{data}$)
 ✓ ?
 Sample new data points \vec{X}

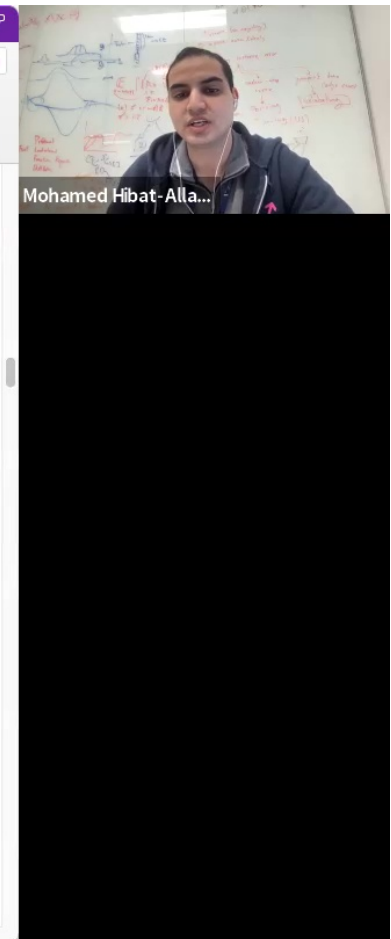
Training data

Generated data

X_2

X_1

MH



OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 0.25 mm 0.35 mm 0.5 mm 0.7 mm 1 mm

$(D_{train} \approx P_{data})$
✓ ?

Sample new data points \vec{x}

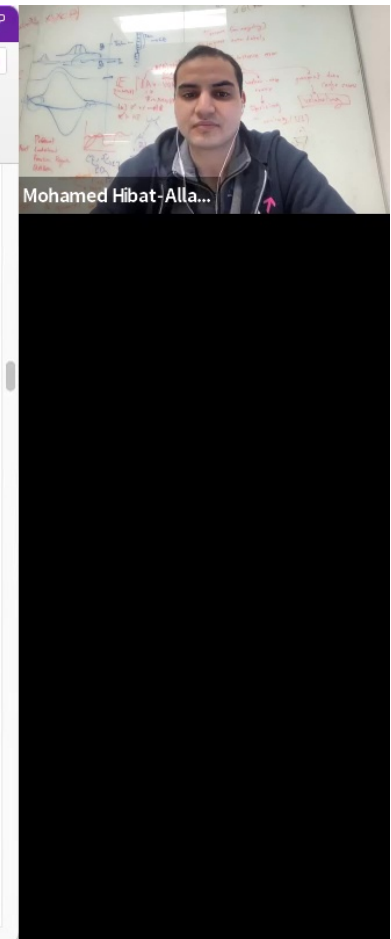
Training data

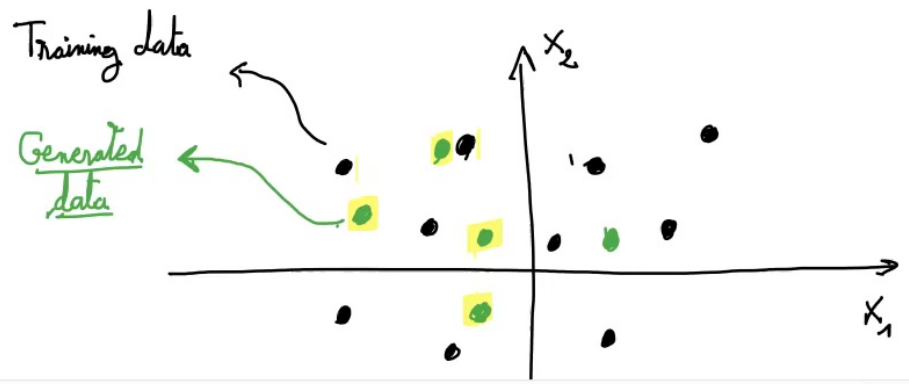
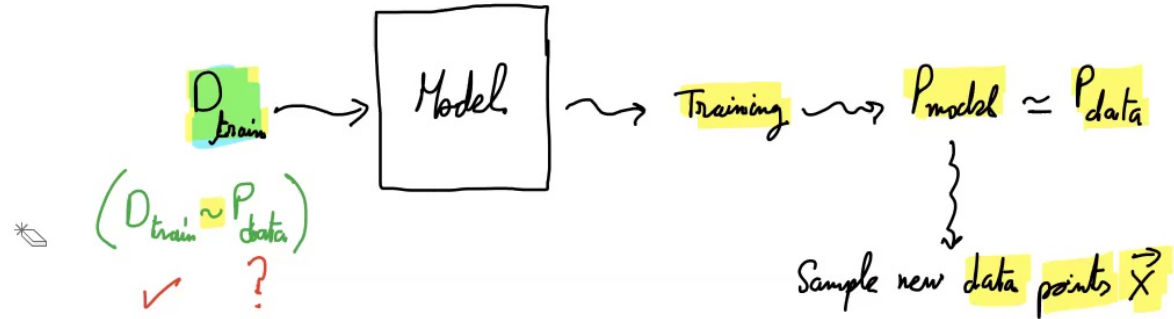
Generated data

x_2

x_1

② Maximum Likelihood Principle:





Perimeter-A


OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 0.25 mm 0.35 mm 0.5 mm 0.7 mm 1 mm

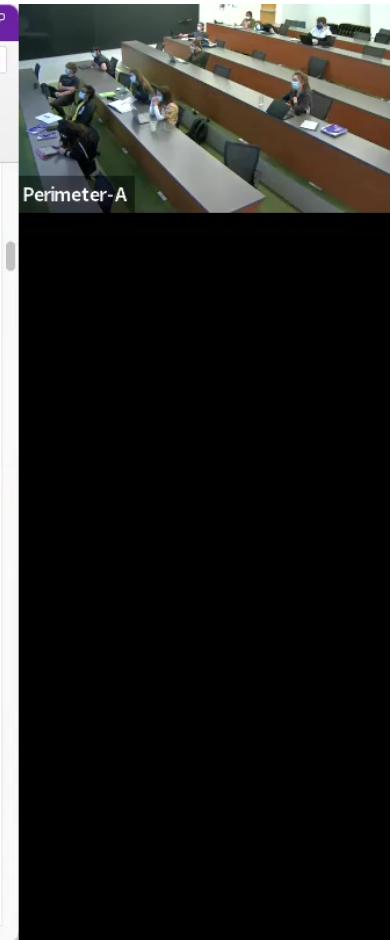
$P_{\text{face}}(\vec{X}) = \dots$

$\vec{X} = \underline{\text{image}}$



Credit: this-person-does-not-exist.com

MH

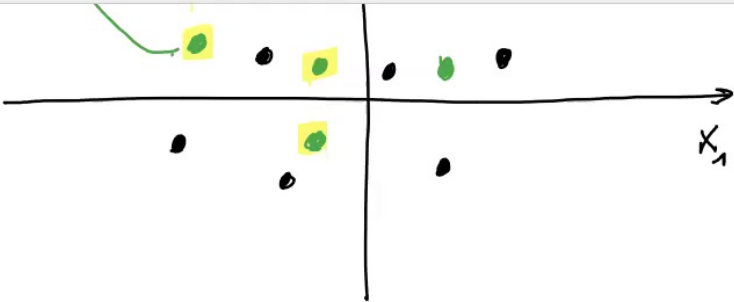


OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 2 mm 3 mm 5 mm 6.5 mm 8 mm

data



② Maximum Likelihood Principle:

↳ Our aim is to find the best model $\vec{\theta}$ given a dataset D

↔ Maximize $P(\vec{\theta} | D)$

↳ Bayes rule:

Likelihood → Prior



② Maximum Likelihood Principle:

→ Our aim is to find the best model $\vec{\theta}$ given a dataset D

↔ Maximize $P(\vec{\theta} | D)$

→ Bayes rule:

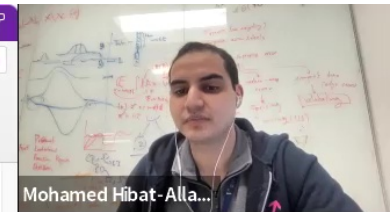
$$P(\vec{\theta} | D) = \frac{P(D | \vec{\theta}) P(\vec{\theta})}{P(D)}$$

Posterior ← Likelihood → Prior
← Evidence →

∝ $P(D | \vec{\theta}) P(\vec{\theta})$



Mohamed Hibat-Alla...

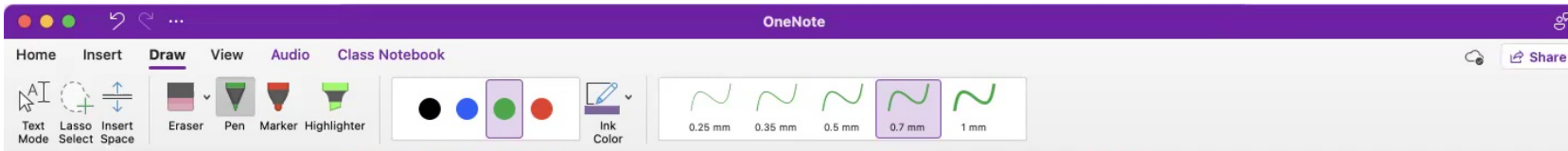


→ Bayes rule:

$$P(\vec{\theta} | D) = \frac{P(D | \vec{\theta}) P(\vec{\theta})}{P(D)}$$

Labels: Posterior (under $P(\vec{\theta} | D)$), Likelihood (under $P(D | \vec{\theta})$), Prior (under $P(\vec{\theta})$), Evidence (under $P(D)$)

- Given a uniform prior $P(\vec{\theta})$
- $\text{Argmax}_{\vec{\theta}} P(\vec{\theta} | D) = \text{Argmax}_{\vec{\theta}} P(D | \vec{\theta})$
- The likelihood $P(D | \vec{\theta})$ is more accessible to calculate.



Bayes rule:

$$P(\vec{\theta} | D) = \frac{P(D | \vec{\theta}) P(\vec{\theta})}{P(D)}$$

Posterior
Evidence
constant

Given a uniform prior $P(\vec{\theta})$

$$\text{Argmax}_{\vec{\theta}} P(\vec{\theta} | D) = \text{Argmax}_{\vec{\theta}} P(D | \vec{\theta})$$

The likelihood $P(D | \vec{\theta})$ is more accessible to calculate.

Example: $P(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$



Mohamed Hibat-Alla...

OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 2 mm 3 mm 5 mm 6.5 mm 8 mm

$$\hookrightarrow \text{Argmax}_{\vec{\theta}} \underline{P(\vec{\theta} | D)} = \text{Argmax}_{\vec{\theta}} \underline{P(D | \vec{\theta})}$$

$$\hookrightarrow \text{The likelihood } P(D | \vec{\theta}) \text{ is more accessible to calculate.}$$

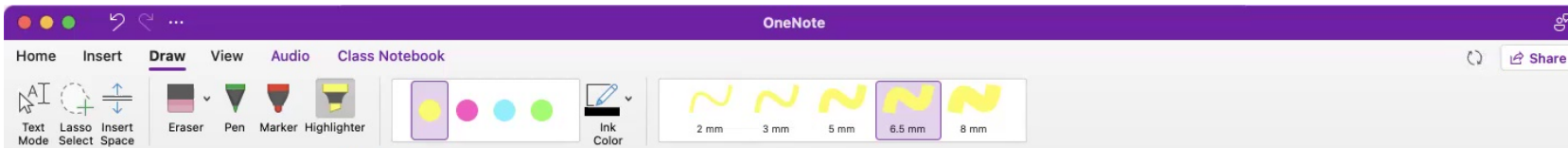
$$\hookrightarrow \text{Example: } P(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(\frac{-1}{2\sigma^2}(x - \mu)^2\right)$$

$$L_{\theta}(D) = P(D | \vec{\theta}) = P(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_M | \vec{\theta})$$

$$D = \left\{ \vec{x}_i \right\}_{i=1}^M \quad \vec{x}_i \text{ are independent} \quad = P(\vec{x}_1 | \theta) P(\vec{x}_2 | \theta) \dots P(\vec{x}_M | \theta) \rightarrow \text{"can be very small"}$$

MH

Mohamed Hibat-Alla...



Mohamed Hibat-Alla...

$D = \left\{ \vec{x}_i \right\}_{i=1}^M$ \vec{x}_i are independent $= P(\vec{x}_1 | \theta) P(\vec{x}_2 | \theta) \dots P(\vec{x}_M | \theta)$ \rightarrow "can be very small"

$\hookrightarrow NLL = -\log(L_{\theta}(D)) = -\sum_{i=1}^M \log(P(\vec{x}_i | \theta))$
 Negative log-likelihood

* Why NLL as a cost function? \rightarrow Kullback-Liebler (KL) divergence

$$KL(P_{data} || P_{\theta}) = \sum_{\vec{x}} P_{data}(\vec{x}) \cdot \log\left(\frac{P_{data}(\vec{x})}{P_{\theta}(\vec{x})}\right)$$



$$\rightarrow NLL = -\log(L_{\theta}(D)) = -\sum_{i=1} \log(P(x_i | \theta))$$

Negative log-likelihood

* Why NLL as a cost function? \rightarrow Kullback-Libler (KL) divergence

$$KL(P_{data} || P_{\theta}) = \sum_{\vec{x}} P_{data}(\vec{x}) \cdot \log\left(\frac{P_{data}(\vec{x})}{P_{\theta}(\vec{x})}\right)$$

$$\left\{ \begin{array}{l} \rightarrow KL(P || q) \geq 0 \end{array} \right.$$

\rightarrow We can use KL as a cost function

$$\left\{ \begin{array}{l} \rightarrow KL(P || q) = 0 \iff P = q \end{array} \right.$$



Mohamed Hibat-Alla...

OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 1.2 mm 1.5 mm 2 mm 3 mm 5 mm

$$KL(P_{data} || P_{\theta}) = \sum_{\vec{x}} P_{data}(\vec{x}) \cdot \log \left(\frac{P_{data}(\vec{x})}{P_{\theta}(\vec{x})} \right)$$

$\rightarrow KL(P || q) \geq 0$
 $\rightarrow KL(P || q) = 0 \iff P = q$

\leadsto We can use KL as a cost function

$$\hookrightarrow KL(P_{data} || P_{\theta}) = \underbrace{\sum_{\vec{x}} P_{data}(\vec{x}) \log(P_{data}(\vec{x}))}_{\text{-entropy of data (constant)}} - \sum_{\vec{x}} P_{data}(\vec{x}) \log(P_{\theta}(\vec{x}))$$

$$KL(P_{data} || P_{\theta}) \approx C + \left[-\frac{1}{M} \sum \log(P_{\theta}(\vec{x})) \right]$$


OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 1.2 mm 1.5 mm 2 mm 3 mm 5 mm

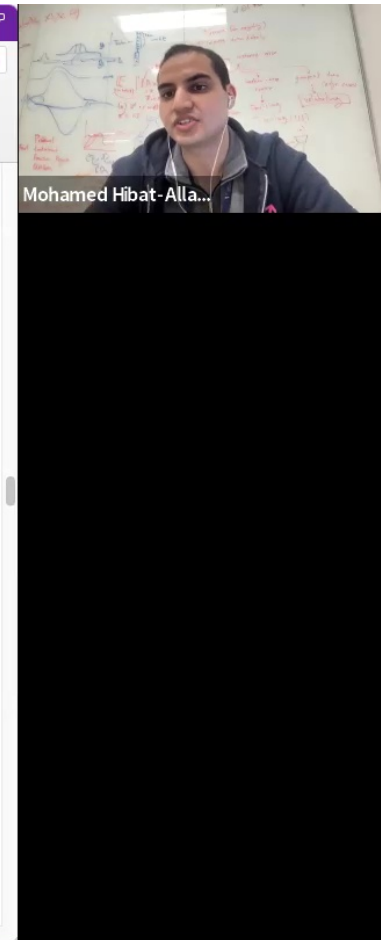
x $\theta(x)$

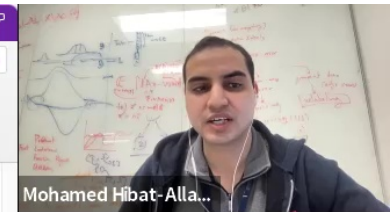
$$\left\{ \begin{array}{l} \rightarrow KL(p \parallel q) \geq 0 \\ \rightarrow KL(p \parallel q) = 0 \Leftrightarrow p = q \end{array} \right.$$

\rightsquigarrow We can use KL as a cost function

$$\hookrightarrow KL(p_{data} \parallel p_{\theta}) = \underbrace{\sum_{\vec{x}} p_{data}(\vec{x}) \log(p_{data}(\vec{x}))}_{- \text{entropy of data (constant)}} - \underbrace{\sum_{\vec{x}} p_{data}(\vec{x}) \log(p_{\theta}(\vec{x}))}_{\approx}$$

$$KL(p_{data} \parallel p_{\theta}) \approx C + \underbrace{\left[-\frac{1}{M} \sum_{i=1}^M \log(p_{\theta}(\vec{x}_i)) \right]}_{\substack{\text{Importance} \\ P(\vec{x}|\theta)}}$$





\vec{x} 'data' $P_{\theta}(\vec{x})$

$$\langle O \rangle = \sum_{\vec{\sigma}} O(\vec{\sigma}) P(\vec{\sigma})$$

$\vec{\sigma} = (\dots)$

$\} \vec{\sigma} \} \sim P$

$$\langle O \rangle \approx \frac{1}{M} \sum_{i=1}^M O(\vec{\sigma}^{(i)})$$

$\rightarrow KL(P||q) \geq 0$
 $\rightarrow KL(P||q) = 0 \iff P = q$

\rightarrow We can use KL as a cost function

$$KL(P_{data} || P_{\theta}) = \underbrace{\sum_{\vec{x}} P_{data}(\vec{x}) \log(P_{data}(\vec{x}))}_{\text{-entropy of data (constant)}} - \underbrace{\sum_{\vec{x}} P_{data}(\vec{x}) \log(P_{\theta}(\vec{x}))}_{\approx}$$

$$KL(P_{data} || P_{\theta}) \approx C + \left[-\frac{1}{M} \sum_{i=1}^M \log(P_{\theta}(\vec{x}^{(i)})) \right]$$

\downarrow Importance $\rightarrow P(\vec{x}|\theta)$

OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Space Eraser Pen Marker Highlighter Ink Color 1.2 mm 1.5 mm 2 mm 3 mm 5 mm

$$KL(P_{data} || P_{\theta}) = \sum_{\vec{x}} P_{data}(\vec{x}) \log(P_{data}(\vec{x})) - \sum_{\vec{x}} P_{data}(\vec{x}) \log(P_{\theta}(\vec{x}))$$

$$\leftarrow \sum_{i=1}^M \frac{1}{M} \log(P_{\theta}(\vec{x}^{(i)}))$$

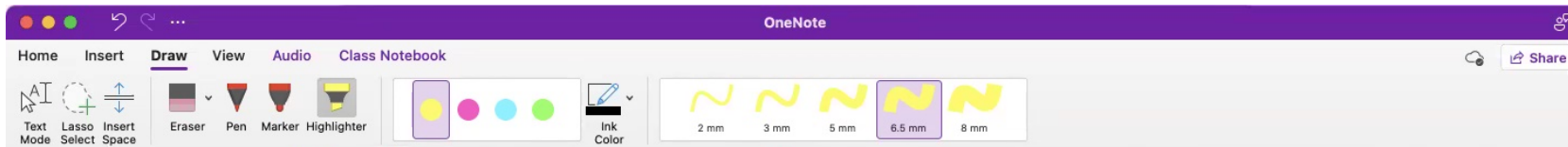
$$KL(P_{data} || P_{\theta}) \approx C + \underbrace{\left[-\frac{1}{M} \sum_{i=1}^M \log(P_{\theta}(\vec{x}^{(i)})) \right]}_{NLL} \rightarrow \underline{P(\vec{x}|\theta)}$$

Importance sampling

$$\rightarrow \text{Minimizing } KL \Leftrightarrow \text{Minimizing } NLL \Leftrightarrow \text{Maximize Likelihood}$$

③ Classification of generative models:





③ Classification of generative models:

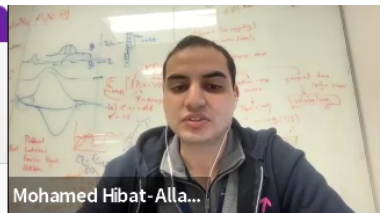
3.1 - Exact likelihood models:

↳ Likelihood $L_{\theta}(\vec{x})$ can be computed exactly.

* Autoregressive models → suitable for discrete data

↳ Probability chain rule:

$$\begin{aligned}
 P_{\theta}(\vec{x}) &= P_{\theta}(x_1) P_{\theta}(x_2|x_1) P_{\theta}(x_3|x_2, x_1) \dots P_{\theta}(x_N|x_{N-2}, \dots, x_2, x_1) \\
 \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow & \\
 P_{\theta}(\vec{x}) &= P_{\theta}(x_1) P_{\theta}(x_2|x_1) P_{\theta}(x_3|x_2, x_1) \dots P_{\theta}(x_N|x_{N-2}, \dots, x_2, x_1)
 \end{aligned}$$



Mohamed Hibat-Alla...

OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 0.25 mm 0.35 mm 0.5 mm 0.7 mm 1 mm

$$P_{\theta}(\vec{x}) = P_{\theta}(x_1) P_{\theta}(x_2|x_1) P_{\theta}(x_3|x_2, x_1) \dots P_{\theta}(x_N|x_{N-1}, \dots, x_2, x_1)$$

$$P_{\theta}(\vec{x}) = P_{\theta}(x_1, x_2, \dots, x_N)$$

$P_{\theta}(\cdot | x_{i-1}, \dots, x_1) = \begin{pmatrix} \vdots \\ \vdots \end{pmatrix}$

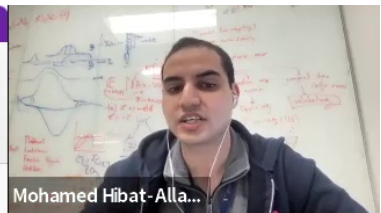
Autoregressive model

Sample x_i

RNN (Next lecture)

Pixel CNNs

Transformers



OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 2 mm 3 mm 5 mm 6.5 mm 8 mm

$$P_{\theta}(\vec{x}) = P_{\theta}(x_1) P_{\theta}(x_2|x_1) P_{\theta}(x_3|x_2, x_1) \dots P_{\theta}(x_N|x_{N-1}, \dots, x_2, x_1)$$

$$P_{\theta}(\vec{x}) = P_{\theta}(x_1, x_2, \dots, x_N)$$

$P_{\theta}(\cdot | x_{i-1}, \dots, x_1) = \begin{pmatrix} \cdot \\ \vdots \end{pmatrix}$

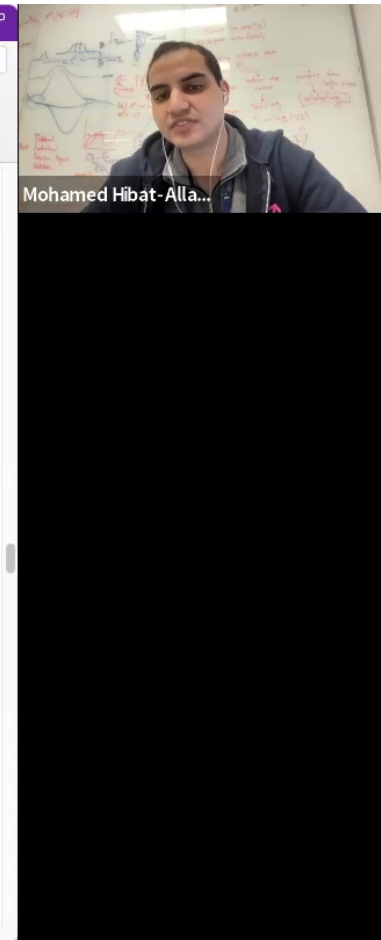
Autoregressive model

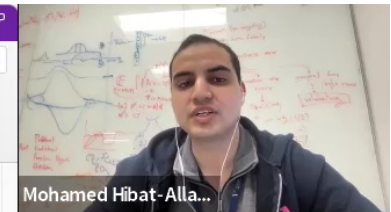
$P(x_1)$

x_0

Sample x_i

RNN (Next lecture)
 Pixel CNNs
 Transformers



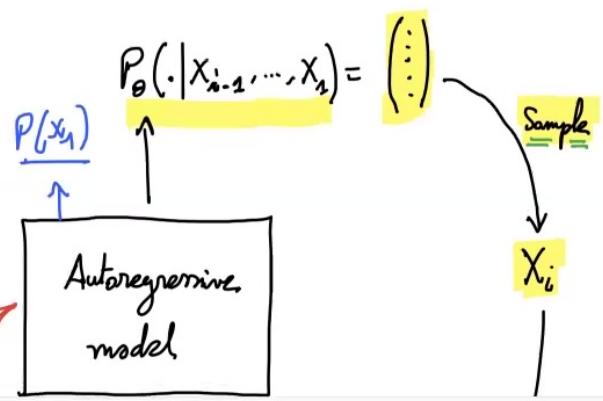


chain

$$P_{\theta}(\vec{x}) = P_{\theta}(x_1) P_{\theta}(x_2|x_1) P_{\theta}(x_3|x_2, x_1) \dots P_{\theta}(x_N|x_{N-2}, \dots, x_2, x_1)$$

$$P_{\theta}(\vec{x}) = P_{\theta}(x_1) P_{\theta}(x_2|x_1) P_{\theta}(x_3|x_2, x_1) \dots P_{\theta}(x_N|x_{N-2}, \dots, x_2, x_1)$$

$$P(\vec{x}) = P_{\theta}(x_1, x_2, \dots, x_N)$$



RNN (Next lecture)

Pixel CNNs

OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 1.2 mm 1.5 mm 2 mm 3 mm 5 mm

Probability chain rule!

$$P_{\theta}(\vec{x}) = P_{\theta}(x_1) P_{\theta}(x_2|x_1) P_{\theta}(x_3|x_2, x_1) \dots P_{\theta}(x_N|x_{N-2}, \dots, x_2, x_1)$$

?

$$P_{\theta}(\vec{x}) = P_{\theta}(x_1) P_{\theta}(x_2|x_1) P_{\theta}(x_3|x_2, x_1) \dots P_{\theta}(x_N|x_{N-2}, \dots, x_2, x_1)$$

$$P_{\theta}(\vec{x}) = P_{\theta}(x_1, x_2, \dots, x_N)$$

Autoregressive model

$P(x_1)$

$P_{\theta}(\cdot|x_{i-2}, \dots, x_1) = \begin{pmatrix} \vdots \\ \vdots \end{pmatrix}$

Sample x_i

RNN (Next lecture)

Pixel CNNs

Transformers



OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color

1.2 mm 1.5 mm 2 mm 3 mm 5 mm

$\rightarrow KL(P||q) = 0 \Leftrightarrow p = q$

$KL(P_{data} || P_{\theta}) = \sum_{\vec{x}} P_{data}(\vec{x}) \log(P_{data}(\vec{x})) - \sum_{\vec{x}} P_{data}(\vec{x}) \log(P_{\theta}(\vec{x}))$

- entropy of data (constant)

$KL(P_{data} || P_{\theta}) \approx C + \left[-\frac{1}{M} \sum_{i=1}^M \log(P_{\theta}(\vec{x}_i)) \right]$

Importance sampling

NLL

$P(\vec{x}|\theta)$

\rightarrow Minimizing KL \Leftrightarrow Minimizing NLL \Leftrightarrow Maximize Likelihood



OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Space Eraser Pen Marker Highlighter Ink Color

1.2 mm 1.5 mm 2 mm 3 mm 5 mm

$\rightarrow KL(P||q) = 0 \Leftrightarrow p = q$

$KL(P_{data} || P_{\theta}) = \sum_{\vec{x}} P_{data}(\vec{x}) \log(P_{data}(\vec{x})) - \sum_{\vec{x}} P_{data}(\vec{x}) \log(P_{\theta}(\vec{x}))$

- entropy of data (constant)

$KL(P_{data} || P_{\theta}) \approx C + \left[-\frac{1}{M} \sum_{i=1}^M \log(P_{\theta}(\vec{x}_i)) \right]$

Importance sampling

$P(\vec{x}|\theta)$

NLL

\rightarrow Minimizing $KL \Leftrightarrow$ Minimizing $NLL \Leftrightarrow$ Maximize Likelihood



OneNote

Home Insert Draw View Audio Class Notebook

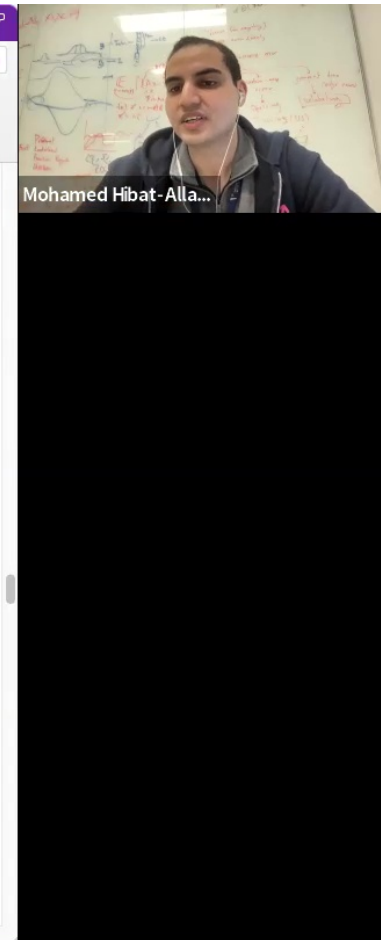
Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 1.2 mm 1.5 mm 2 mm 3 mm 5 mm

* Size of $P(\cdot | x_{c,i})$ depends on the task:

- Words: 1K - 1M
- Pixels: 256×3 (RGB)
- Spins: 2 (\downarrow or \uparrow)

* Normalizing Flows \rightsquigarrow Suitable for continuous data

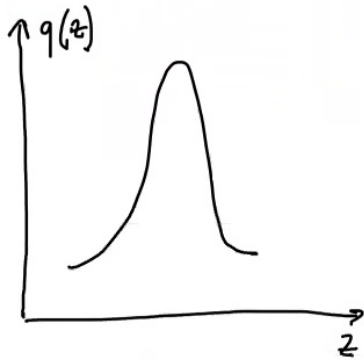
The diagram illustrates the concept of normalizing flows. On the left, a smooth curve is labeled $q(z)$. An arrow labeled "Coordinate" points to the right, where a jagged curve is labeled $P(x)$. This represents the transformation of a simple latent distribution into a complex data distribution.



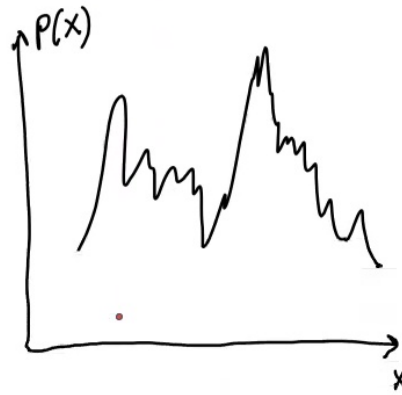


→ Spino: $\mathcal{L}(\downarrow \uparrow)$

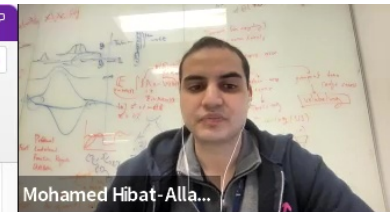
* Normalizing Flows \rightsquigarrow Suitable for continuous data



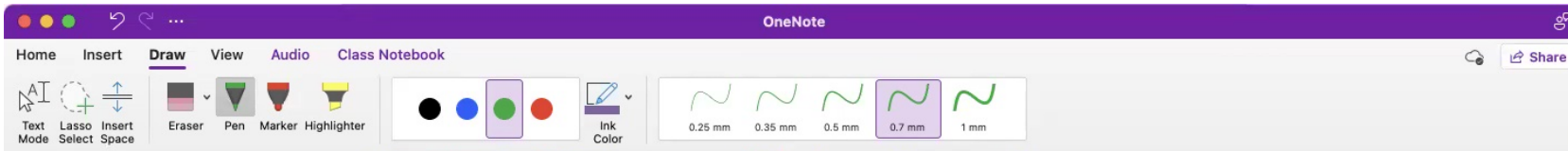
Coordinate
Transformation
 $x = g(z)$



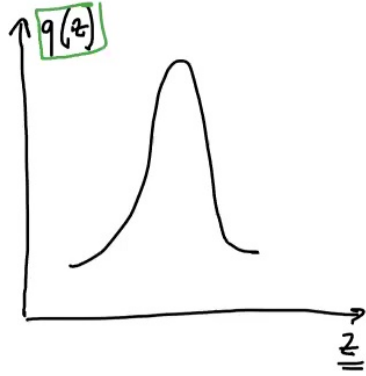
$$\int q(z) dz = 1 \rightsquigarrow \int q(g^{-1}(x)) \cdot \left| \frac{dz}{dx} \right| dx = 1$$



Mohamed Hibat-Alla...



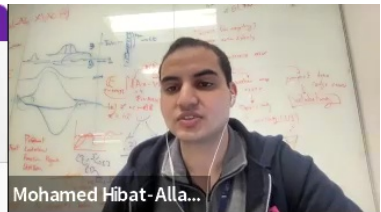
* Normalizing Flows \rightsquigarrow Suitable for continuous data



Coordinate
Transformation
 $X = g(z)$



$$\int q(z) dz = 1 \rightsquigarrow \int \underbrace{q(g^{-1}(x)) \cdot \left| \frac{\partial z}{\partial x} \right|}_{P(x)} dx = 1$$



Mohamed Hibat-Alla...

OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 0.25 mm 0.35 mm 0.5 mm 0.7 mm 1 mm

Coordinate Transformation

$x = g(z)$

$z = g^{-1}(x)$ bijective

$\int q(z) dz = 1 \rightsquigarrow \int \underbrace{q(g^{-1}(x)) \cdot \left| \frac{dz}{dx} \right|}_{p(x)} dx = 1$

$\rightsquigarrow p(x) = q(z) \left| \frac{dz}{dx} \right|$



OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 1.2 mm 1.5 mm 2 mm 3 mm 5 mm

$$P(x) = q(z) \left| \frac{\partial z}{\partial x} \right|$$

* For multi-variables:
$$P(\vec{x}) = q(\vec{z}) \left| \det \left(\frac{\partial \vec{z}}{\partial \vec{x}} \right) \right|$$
 Jacobian

* The function q is chosen carefully, such that " $\det \left(\frac{\partial g^{-1}(\vec{x})}{\partial \vec{x}} \right)$ " can be computed efficiently

$$\frac{\partial g^{-1}(\vec{x})}{\partial \vec{x}} = \begin{pmatrix} * & & & \\ & * & & \\ & & \dots & \\ * & & & * \end{pmatrix}$$

MH



OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 1.2 mm 1.5 mm 2 mm 3 mm 5 mm

\vec{x}

\ast

\ast

Example:

$x_1 = z_1$

$x_2 = z_2 \exp(s(z_1)) + t(z_1)$

$z_1 = x_1$

$z_2 = (x_2 - t(x_1)) \cdot \exp(-s(x_1))$

$\rightarrow \det \left(\frac{\partial \vec{z}}{\partial \vec{x}} \right) = \det \begin{pmatrix} 1 & 0 \\ 0 & \exp(-s(x_1)) \end{pmatrix} = \exp(-s(x_1))$



OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 2 mm 3 mm 5 mm 6.5 mm 8 mm

$$x_2 = z_2 \exp(s(z_1)) + t(z_1)$$

$$z_2 = x_2 - t(x_1)$$

$$z_2 = (x_2 - t(x_1)) \cdot \exp(-s(x_1))$$

$$\det \left(\frac{\partial \vec{z}}{\partial \vec{x}} \right) = \det \begin{pmatrix} 1 & 0 \\ 0 & \exp(-s(x_1)) \end{pmatrix} = \underline{\underline{\exp(-s(x_1))}}$$

→ This argument can be generalized for N variables.

→ This transformation can be iterated $\vec{x} = \underbrace{f_1 \circ f_2 \circ \dots \circ f_L}_{\mathcal{g}}(\vec{z})$



OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 1.2 mm 1.5 mm 2 mm 3 mm 5 mm

→ Words: 1K - 1M
 → Pixels: 256 x 3 (RGB)
 → Spins: 2 (↓ or ↑)

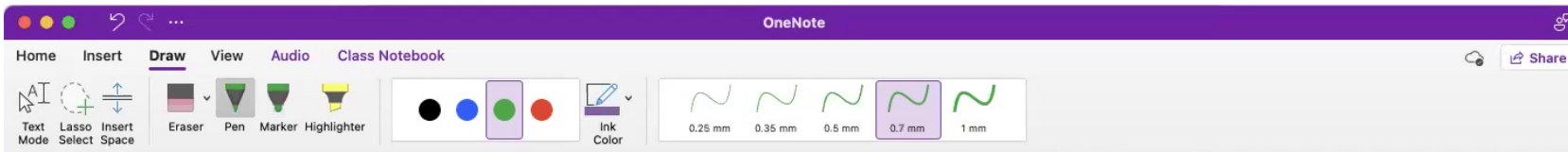
* Normalizing Flows → Suitable for continuous data

$q(z)$

Coordinate Transformation
 $x = g(z)$

$P(x) \approx P_{data}$





→ This transformation can be iterated $\vec{x} = \underbrace{f_1 \circ f_2 \circ \dots \circ f_L}_{g}(\vec{z})$

* 3.2 - Approximate likelihood models:

↳ $L_{\theta}(\vec{x})$ is intractable to compute.

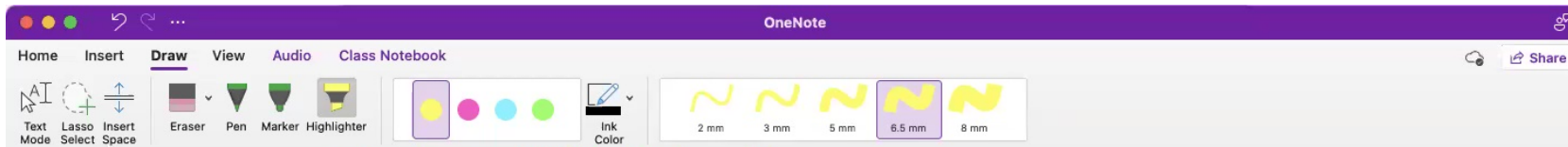
* Energy-based models → Restricted Boltzmann Machines (RBM)

$$\hookrightarrow E(\vec{x}, \vec{h}) = \sum_{i=1}^N \sum_{j=1}^M x_i w_{ij} h_j - \sum_{i=1}^N a_i x_i - \sum_{j=1}^M b_j h_j$$

→ l l l h h_s



Mohamed Hibat-Alla...

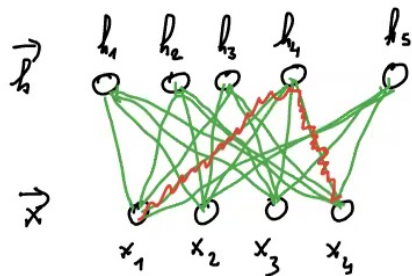


* 3.2 Approximate likelihood models:

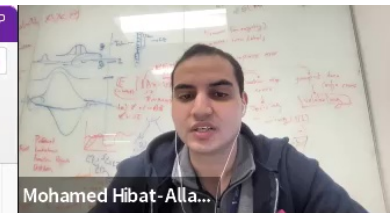
↳ $\mathcal{L}(\vec{X})$ is intractable to compute.

* Energy-based models → Restricted Boltzmann Machines (RBM)

$$E(\vec{X}, \vec{h}) = \sum_{i=1}^N \sum_{j=1}^M x_i w_{ij} h_j - \sum_{i=1}^N a_i x_i - \sum_{j=1}^M b_j h_j$$



→ Correlation between x_1 and x_4 is introduced through h_4 .

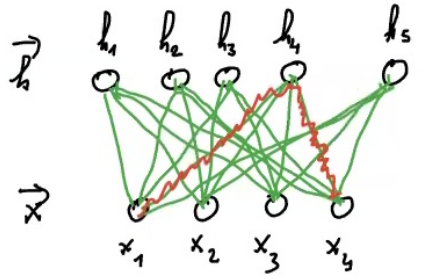


Mohamed Hibat-Alla...

* Energy-based models \rightarrow Restricted Boltzmann Machines (RBM)

$\vec{x} = () \uparrow N$
 $\vec{h} = () \uparrow M$

$$E(\vec{x}, \vec{h}) = \sum_{i=1}^N \sum_{j=1}^M x_i w_{ij} h_j - \sum_{i=1}^N a_i x_i - \sum_{j=1}^M b_j h_j$$

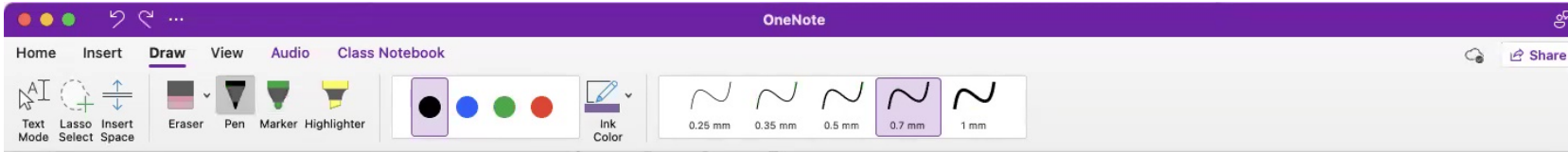


\rightarrow Correlation between x_1 and x_4 is introduced through h_4 .

$$\rho(\vec{x}, \vec{h}) = \frac{\exp(-E(\vec{x}, \vec{h}))}{Z} ; Z = \sum_{\vec{x}, \vec{h}} \exp(-E(\vec{x}, \vec{h}))$$



Mohamed Hibat-Alla...



Mohamed Hibat-Alla...

$$\leadsto p(\vec{x}, \vec{h}) = \frac{\exp(-E(\vec{x}, \vec{h}))}{Z} ; Z = \sum_{\vec{x}, \vec{h}} \exp(-E(\vec{x}, \vec{h}))$$

$$\leadsto P(x) = \left[\sum_{\vec{h}} \exp(-E(\vec{x}, \vec{h})) \right] \times \frac{1}{Z}$$

Numerator

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \sum_{h_1=\pm 1} \dots \sum_{h_N=\pm 1} \exp\left(-\left(\sum_i x_i w_{i1} + b_1\right) h_1\right) \times \dots \times \exp\left(-\left(\sum_i x_i w_{iN} + b_N\right) h_N\right)$$

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \prod_{j=1}^M \left(\sum_{h_j=\pm 1} \exp\left(-\left(\sum_i x_i w_{ij} + b_j\right) h_j\right) \right)$$

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \prod_{j=1}^M \left(2 \cosh\left(\sum_i x_i w_{ij} + b_j\right) \right)$$

OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 2 mm 3 mm 5 mm 6.5 mm 8 mm

Numerator

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \sum_{h_1=\pm 1} \dots \sum_{h_M=\pm 1} \underbrace{\exp\left(-\left(\sum_i x_i w_{i1} + b_1\right) h_1\right)}_{\leftarrow} \dots \underbrace{\exp\left(-\left(\sum_i x_i w_{iN} + b_N\right) h_N\right)}_{\leftarrow}$$

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \prod_{j=1}^M \left(\sum_{h_j=\pm 1} \exp\left(-\left(\sum_i x_i w_{ij} + b_j\right) h_j\right) \right)$$

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \prod_{j=1}^M \left(2 \cosh\left(\sum_i x_i w_{ij} + b_j\right) \right)$$

$Z = \sum_i (\text{Numerator}) \rightsquigarrow$ Intractable \rightsquigarrow Contrastive Divergence (Approximation)

$\hookrightarrow P(\vec{x})$ can be sampled using Metropolis sampling



OneNote

Home Insert Draw View Audio Class Notebook Share

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 2 mm 3 mm 5 mm 6.5 mm 8 mm

Numerator

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \sum_{h_1=\pm 1} \dots \sum_{h_M=\pm 1} \exp\left(-\left(\sum_i x_i w_{i1} + b_1\right) h_1\right) \dots \exp\left(\left(\sum_i x_i w_{iN} + b_N\right) h_N\right)$$

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \prod_{j=1}^M \left(\sum_{h_j=\pm 1} \exp\left(-\left(\sum_i x_i w_{ij} + b_j\right) h_j\right) \right)$$

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \prod_{j=1}^M \left(2 \cosh\left(\sum_i x_i w_{ij} + b_j\right) \right)$$

$Z = \sum_i (\text{Numerator}) \rightsquigarrow$ Intractable \rightsquigarrow Contrastive Divergence (Approximation)

$\hookrightarrow P(\vec{x})$ can be sampled using Metropolis sampling

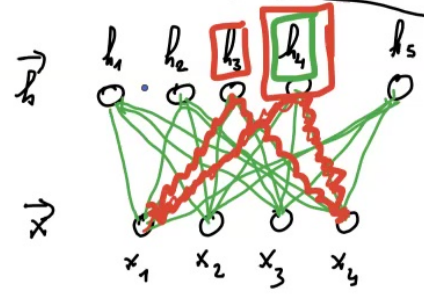


* Energy-based models \rightarrow Restricted Boltzmann Machines (RBM)

$$E(\vec{x}, \vec{h}) = \sum_{i=1}^N \sum_{j=1}^M x_i w_{ij} h_j - \sum_{i=1}^N a_i x_i - \sum_{j=1}^M b_j h_j$$

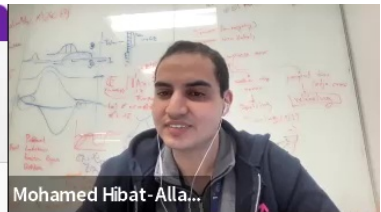
$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} \uparrow N$

$\vec{h} = \begin{pmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \\ h_5 \end{pmatrix} \uparrow M$



\rightarrow Correlation between x_1 and x_4 is introduced through h_4 .

$$p(\vec{x}, \vec{h}) = \frac{\exp(-E(\vec{x}, \vec{h}))}{Z} ; Z = \sum_{\vec{x}, \vec{h}} \exp(-E(\vec{x}, \vec{h}))$$



Mohamed Hibat-Alla...

OneNote

Home Insert Draw View Audio Class Notebook

Text Mode Lasso Select Insert Space Eraser Pen Marker Highlighter Ink Color 1.2 mm 1.5 mm 2 mm 3 mm 5 mm

$$\hookrightarrow \text{Loss: } L = \|\hat{X} - X\|^2$$

- * By take new values of Z , we can generate new data points using Decoder
- * Reducing X to Z allow to compress data and to determine relevant features.
- * The probabilistic interpretation is omitted to keep the discussion simple (See Variational Autoencoders)

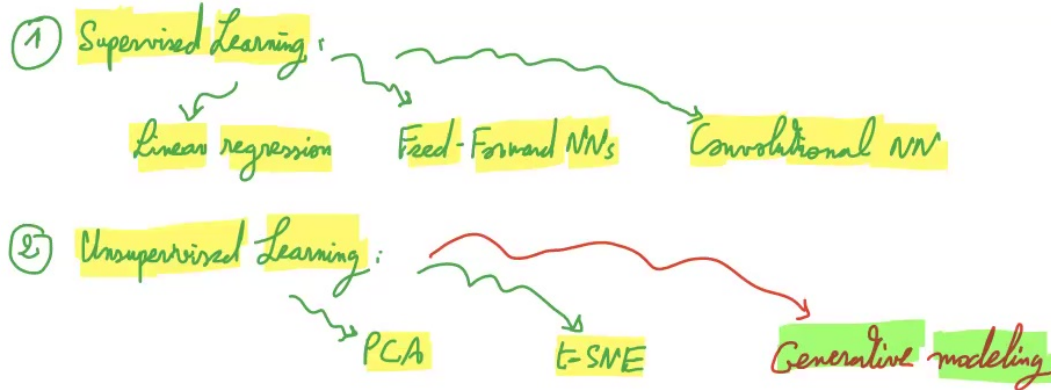




Introduction to generative modeling

Introduction to generative modeling

↳ Review :



↳ Outline :



$$\leadsto \underline{p(\vec{x}, \vec{h})} = \frac{\exp(-E(\vec{x}, \vec{h}))}{Z} ; Z = \sum_{\vec{x}, \vec{h}} \exp(-E(\vec{x}, \vec{h}))$$

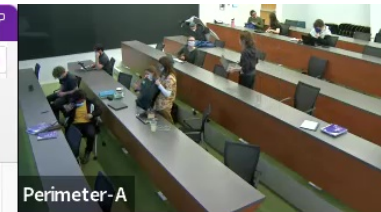
$$\leadsto P(x) = \left[\sum_{\vec{h}} \exp(-E(\vec{x}, \vec{h})) \right] \times \frac{1}{Z}$$

Numerator

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \sum_{h_2=\pm 1} \dots \sum_{h_M=\pm 1} \exp\left(-\left(\sum_i x_i w_{i2} + b_2\right) h_2\right) \dots \exp\left(-\left(\sum_i x_i w_{iN} + b_N\right) h_N\right)$$

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \prod_{j=1}^M \left(\sum_{h_j=\pm 1} \exp\left(-\left(\sum_i x_i w_{ij} + b_j\right) h_j\right) \right)$$

$$= \exp\left(-\sum_{i=1}^N a_i x_i\right) \prod_{j=1}^M \left(2 \cosh\left(\sum_i x_i w_{ij} + b_j\right) \right)$$



Perimeter-A