

Title: Restricted Boltzmann Machines (RBMs)

Speakers: Mohamed Hibat Allah

Collection: Navigating Quantum and AI Career Trajectories: A Beginner's Mini-Course on Computational Methods and their Applications

Date: May 21, 2024 - 11:30 AM

URL: <https://pirsa.org/24050083>

Lecture 1

Restricted Boltzmann Machines (RBMs)

* Outline:

- ① Motivation for Generative modeling.
- ② Generative models classification.
- ③ Restricted Boltzmann Machines (RBMs).

2017-1

Restricted Boltzmann Machines (RBMs)

* Outline:

- ① Motivation for Generative modeling.
- ② Generative models classification.
- ③ Restricted Boltzmann Machines (RBMs).

↳ Machine Learning references:

- ② Generative models classification.
- ③ Restricted Boltzmann Machines (RBMs).

↳ Machine Learning references:

- ↳ Machine Learning for many-body physics, PIRSA course, pirsa.org/c24027
- ↳ "Deep learning specialization", Coursera, Andrew Ng
- ↳ "Neural Networks and Deep Learning", Nielson

* uniqueness for this version.

↳ "A Practical guide to Training Restricted Boltzmann Machines". G. Hinton

↳ "A high-bias, low-variance introduction to machine learning for physicists", ArXiv:1803.08823

↳ "Restricted Boltzmann Machines: A short Tutorial" by Perimeter Institute Quantum Intelligence

https://qucumber.readthedocs.io/en/stable/_static/RBM_tutorial.pdf

① Motivation for generative modeling: → very relevant applications to QPhys

↳ It's a machine learning strategy among many strategies

→ Supervised learning $D = \{(\vec{x}, \vec{y})\}$
Data points Labels

↳ Learn f such that $f(\vec{x}) \approx \vec{y}$

↳ Infer $f(\vec{x}_{\text{new}})$ New datapoints

→ Supervised learning

$$D = \{(\vec{x}, \vec{y})\}$$

Data points Labels

↳ Learn f such that $f(\vec{x}) \approx \vec{y}$

↳ Infer $f(\vec{x}_{\text{new}})$ New datapoints
← Good generalization

→ Unsupervised learning:

$$D = \{\vec{x}\} \leftarrow \text{No labels}$$

Dimensionality reduction

Generative modeling

Learn f such that $f(\vec{x}) \approx \vec{y}$

Infer $f(\vec{x}_{new})$

data points (pointing to \vec{x})

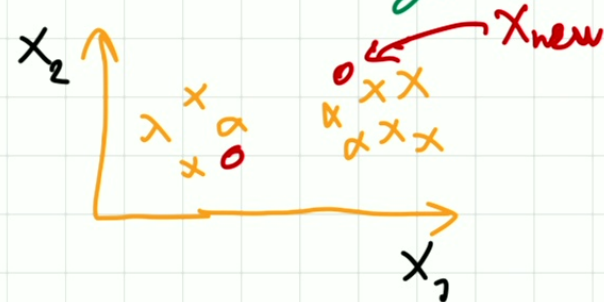
New datapoints (pointing to \vec{x}_{new})

Good generalization (pointing to the inference step)

→ Unsupervised learning: $D = \{ \vec{x} \}$ ← *No labels*

Dimensionality reduction

Generative modeling



→ Reinforcement learning:

↳ Infer $f(\vec{x}_{new})$ ← New datapoints

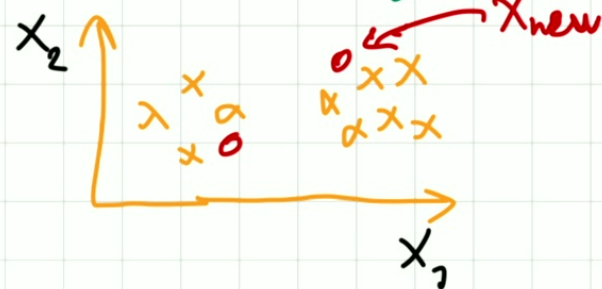
← Good generalization

→ Unsupervised learning:

$D = \{ \vec{x} \}$ ← No labels

↳ Dimensionality reduction

↳ Generative modeling



→ Reinforcement learning:

↳ We aim to use a ML model to maximize
a reward r

* Generative modeling:

e.g. GAN (deep fake), diffusion models (DALL-E)

ChatGPT, GPT-4 → cf. Ethics session at 2pm

* Principle:

↳ We have a dataset $\{\vec{x}\}$, what's $P_{\text{data}}(\vec{x})$?
sampled from $P_{\text{data}}(\vec{x})$

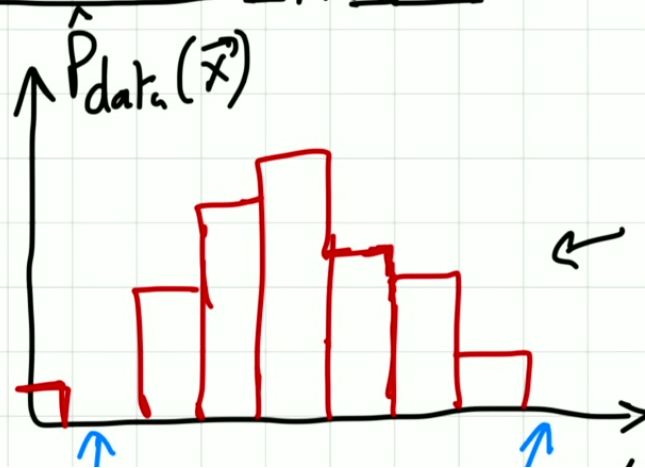
- Traditional approach: (statistics)

Chat GPT, GPT-4 → cf. Ethics session at 2pm

* Principle:

↳ We have a dataset $\{\vec{x}\}$, what's $P_{\text{data}}(\vec{x})$?
sampled from $P_{\text{data}}(\vec{x})$

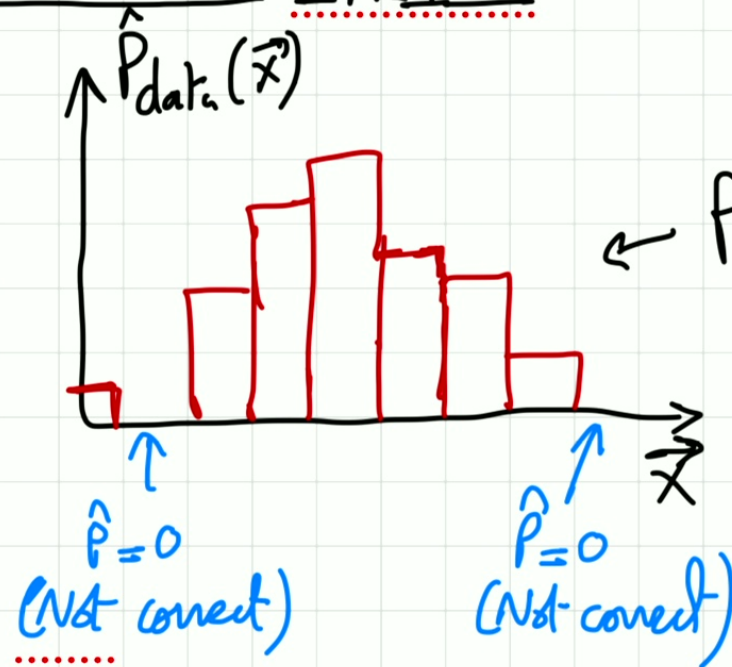
→ Traditional approach: (Statistics)



$$P(\vec{x}) \approx \frac{1}{|D|} \sum_{\vec{x} \in D} \delta(\vec{x} - \vec{x}_k)$$

↳ We have a dataset $\{\vec{x}^i\}_i$, what's $P_{\text{data}}(\vec{x})$?

→ Traditional approach: (Statistics)

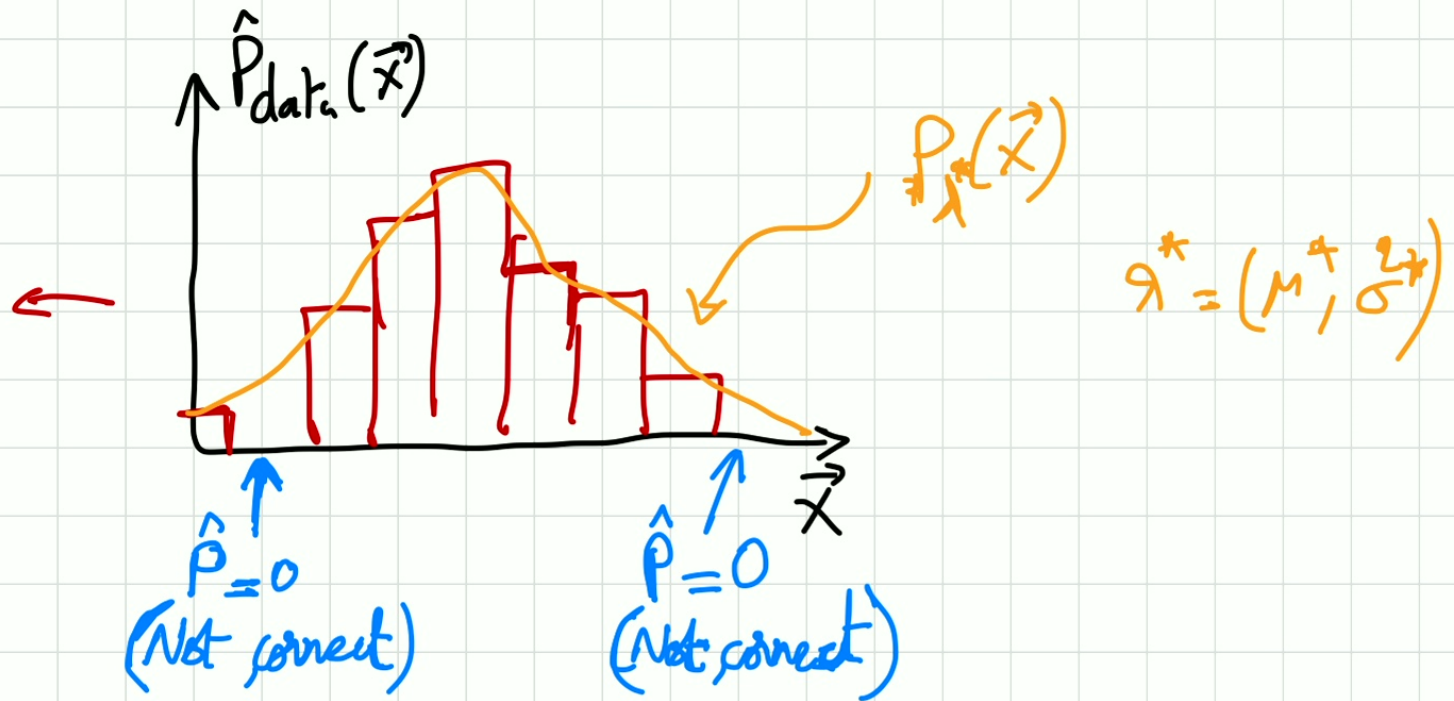


$$P(\vec{x}) \approx \frac{1}{|\mathcal{D}|} \sum_{\vec{x}^i \in \mathcal{D}} \delta(\vec{x} - \vec{x}^i)$$

Two parameters

|| ||

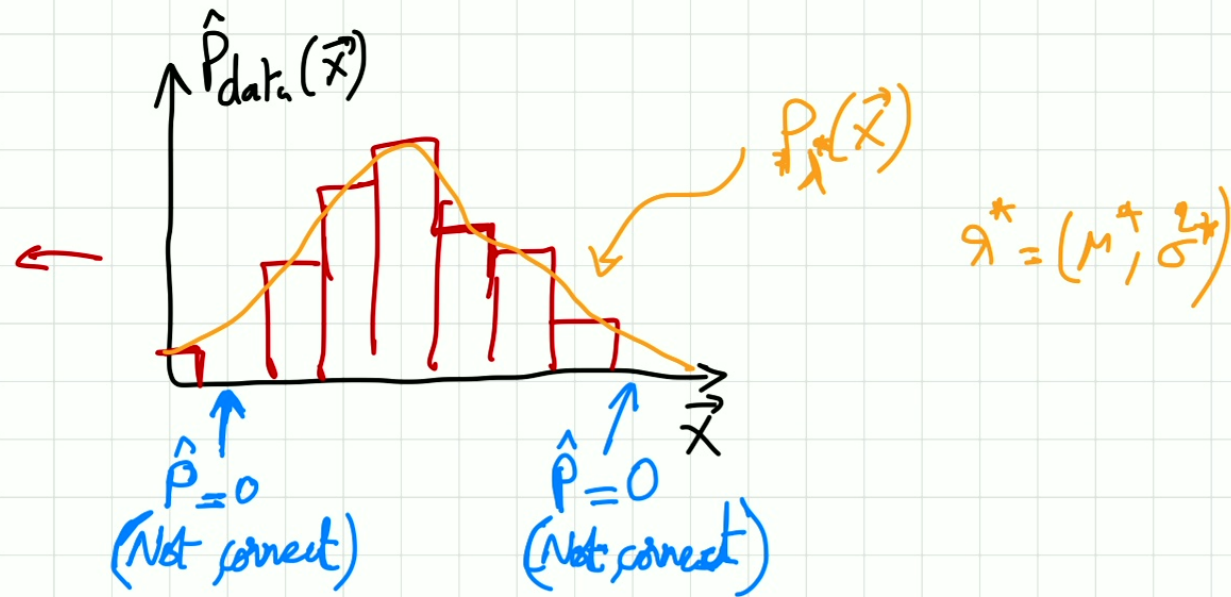
See demo
on Colab
.....



→ μ^*, σ^* is obtained using the maximum likelihood principle

$$\vec{\theta}^* = \operatorname{argmax}_{\vec{\theta}} P_{\vec{\theta}}(D)$$

See demo
on Colab
.....



→ μ^*, σ^* is obtained using the maximum likelihood principle

$$\vec{\lambda}^* = \underset{\vec{\lambda}}{\operatorname{argmax}} \underbrace{P_{\vec{\lambda}}(D)}_{\prod_{i=1}^M P_{\vec{\lambda}}(\vec{X}_i)}$$

Assuming
independence

MLE_principle_gaussian_distribution.ipynb

File Edit View Insert Runtime Tools Help [All changes saved](#)

Comment Share

+ Code + Text

Connect T4

↑ ↓ ↻ ⌨ ⚙ 📄 🗑 ⋮

```
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns |
from scipy.stats import norm

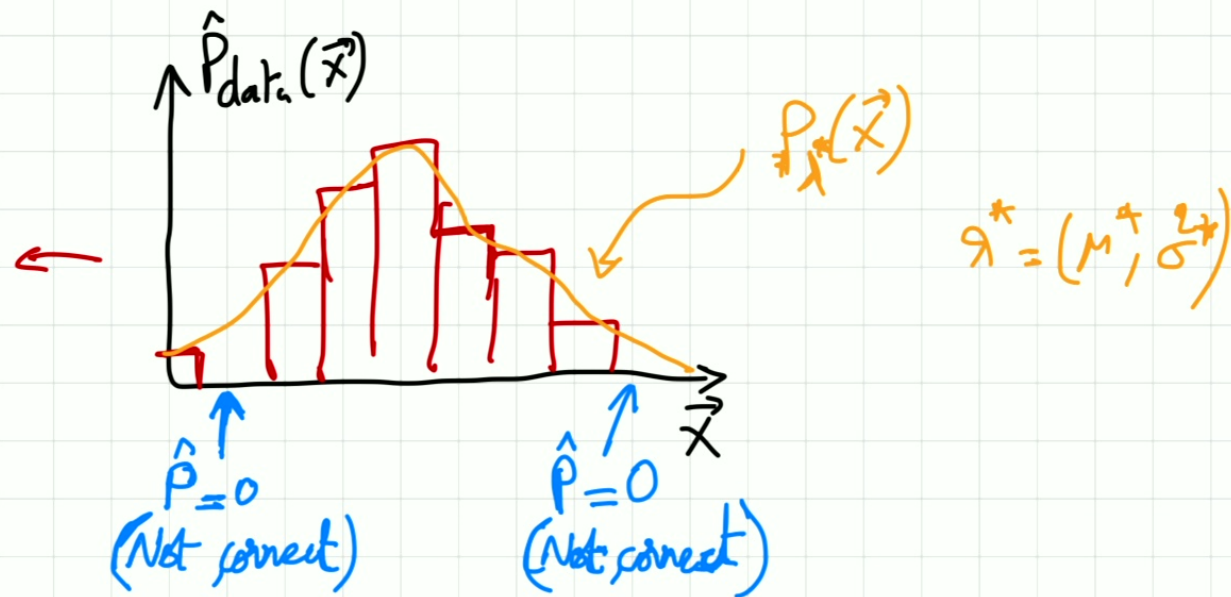
# True parameters
mu_true = 1
sigma_true = 1

# Generate random data
np.random.seed(0)
data = np.random.normal(mu_true, sigma_true, 100)

[ ] initial_guess = [0.5, 0.1] #Initial ansatz for mu and sigma

[ ] sns.histplot(data, kde=False, stat="density", label='Histogram')
x = np.linspace(-2, 4, 300)
plt.plot(x, norm.pdf(x, mu_true, sigma_true), 'r-', lw=2, label='True mu, sigma')
```

See demo
on Colab



→ μ^*, σ^* is obtained using the maximum likelihood principle

$$\vec{\lambda}^* = \underset{\vec{\lambda}}{\operatorname{argmax}} \underbrace{P_{\vec{\lambda}}(D)}_{\prod_{i=1}^M P_{\vec{\lambda}}(\vec{x}_i)}$$

Assuming
|| ||

$P=0$
(Not correct)

$P=0$
(Not correct)

→ μ^*, σ^* is obtained using the maximum likelihood principle

$$\vec{\lambda}^* = \underset{\vec{\lambda}}{\operatorname{argmax}} \underbrace{P_{\vec{\lambda}}(D)}_{\prod_{i=1}^M P_{\vec{\lambda}}(\vec{X}_i)}$$

Assuming
datapoints are
independent

* Why do we do this?

↳ Efficient representation of $p(\vec{x})_{\text{data}}$



$$\prod_{i=1}^M P_{\vec{x}}(\vec{x}_i)$$

Assuming
datapoints are
independent

* Why do we do this?

↳ Efficient representation of $p(\vec{x})_{\text{data}}$

↳ Calculating properties / Observables $\langle O \rangle_{p(\vec{x})_{\text{data}}}$

↳ Ability to generalize to unseen data.

⊗ In quantum physics:

↳ Adapt generative models to construct $\psi_{\lambda}(\vec{x}) = \sqrt{P_{\lambda}(\vec{x})}$
 $\propto \exp(i\phi_{\lambda}(\vec{x}))$

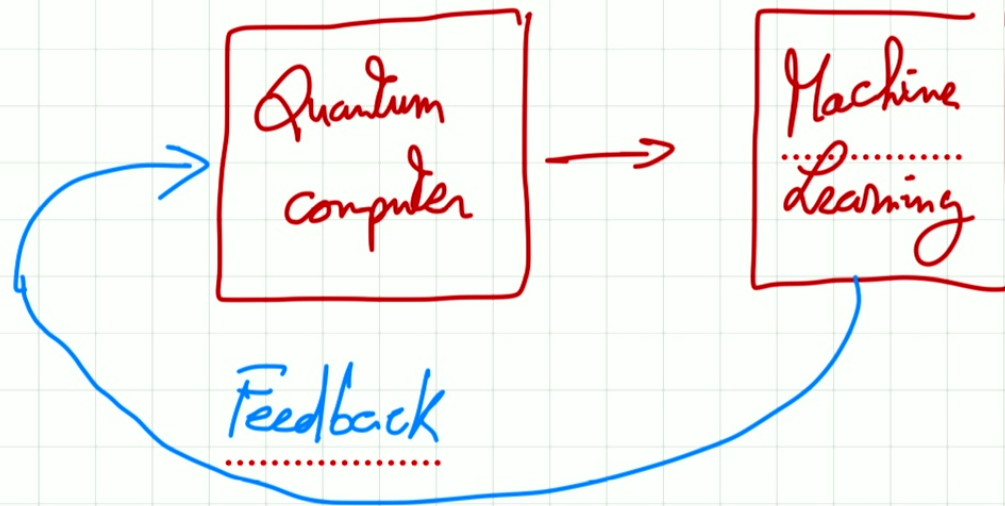
* Goal: find an ML model in $D > 1$ better than SOTA.

- Tensor Networks are very effective in $D=1$.

→ Tensor Networks are very effective in $D=1$.

Another Goal:

↳ Advances the development of Quantum Computers.



② Classification of generative models:

Explicit density

Direct parametrization

Tractable

« Recurrent
Neural
Networks
(RNNs) »

→ Thurs

« Transformers »

→ Friday

Intractable

« Restricted Boltzmann
Machines (RBMs) »

↓
Today

Implicit density

No direct parametrization

(Generative Adversarial
Network)

③ Restricted Boltzmann machines (RBM)

↳ Belongs to the class of energy-based models

$$P_{\vec{\lambda}}(\vec{x}) = \frac{\exp(-E_{\vec{\lambda}}(\vec{x}))}{Z_{\vec{\lambda}}} \leftarrow \text{Partition function}$$

$$\text{where } Z_{\vec{\lambda}} = \sum_{\vec{x}} \exp(-E_{\vec{\lambda}}(\vec{x}))$$

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

\swarrow
 $0, \pm 1$

Z_{λ} ← Partition function

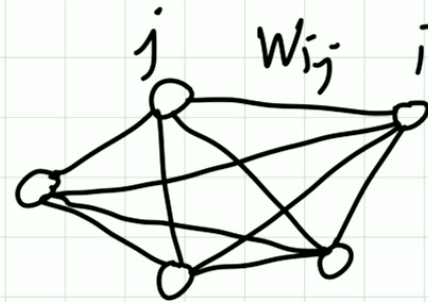
where $Z_{\lambda} = \sum_{\vec{x}} \exp(-E_{\lambda}(\vec{x}))$

$\underbrace{\vec{x}}_{2^N}$

↳ J. Hopfield (1982):

$$E_{\lambda}(\vec{x}) = \sum_{ij} W_{ij} x_i x_j$$

Parameters $\in \mathbb{R}$



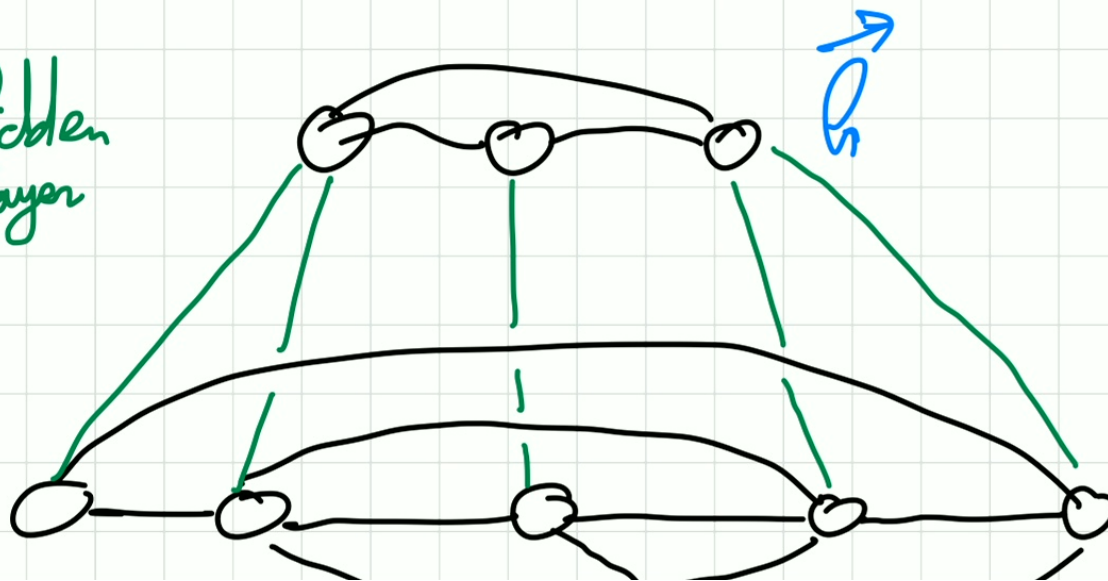
* Disadvantages:

This motivates →

Boltzmann machine (1985): Hinton, Ackley,
"A learning algorithm for Boltzmann
Machines" Sejnowski

could also be
↑ ±1

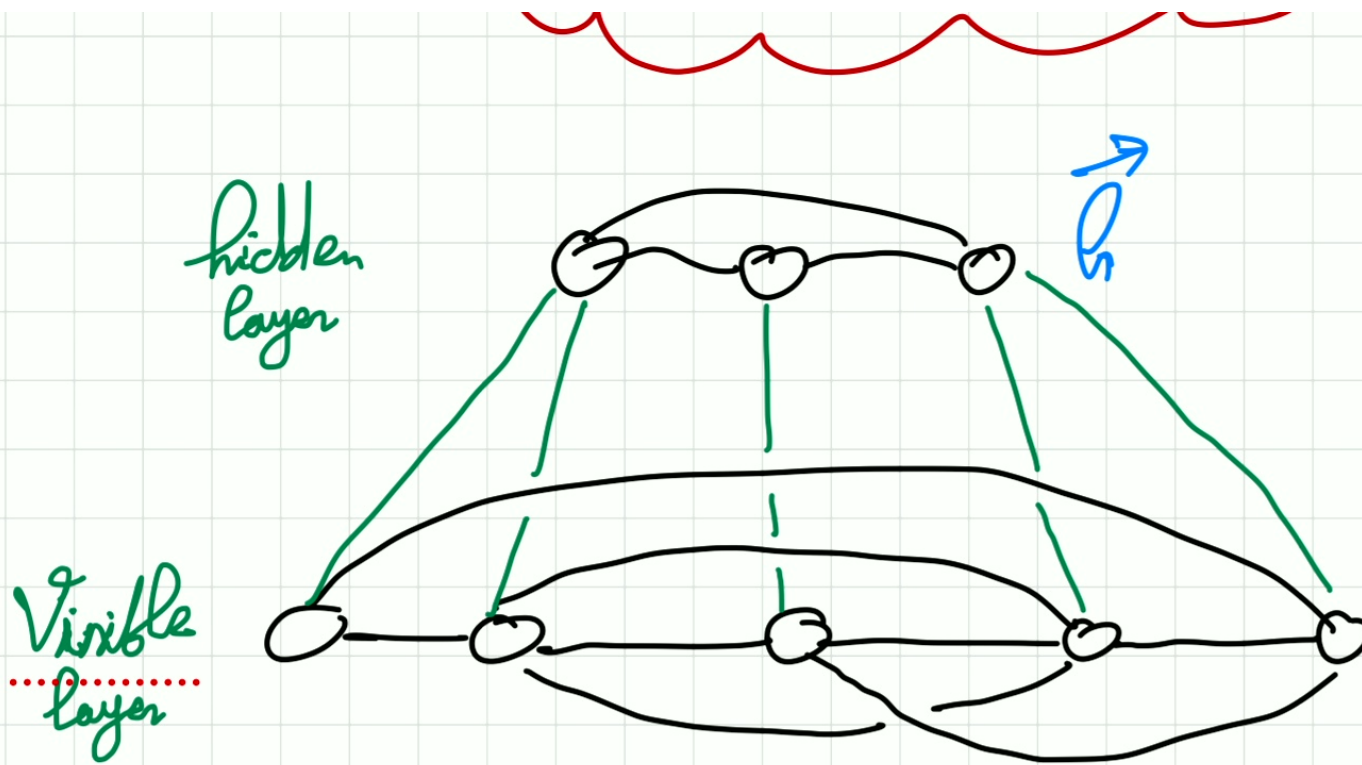
hidden
layer



$h_i = 0, 1$

$x_i = 0, 1$

n



could also be

± 1

$h_i = 0, 1$

$x_i = 0, 1$

n

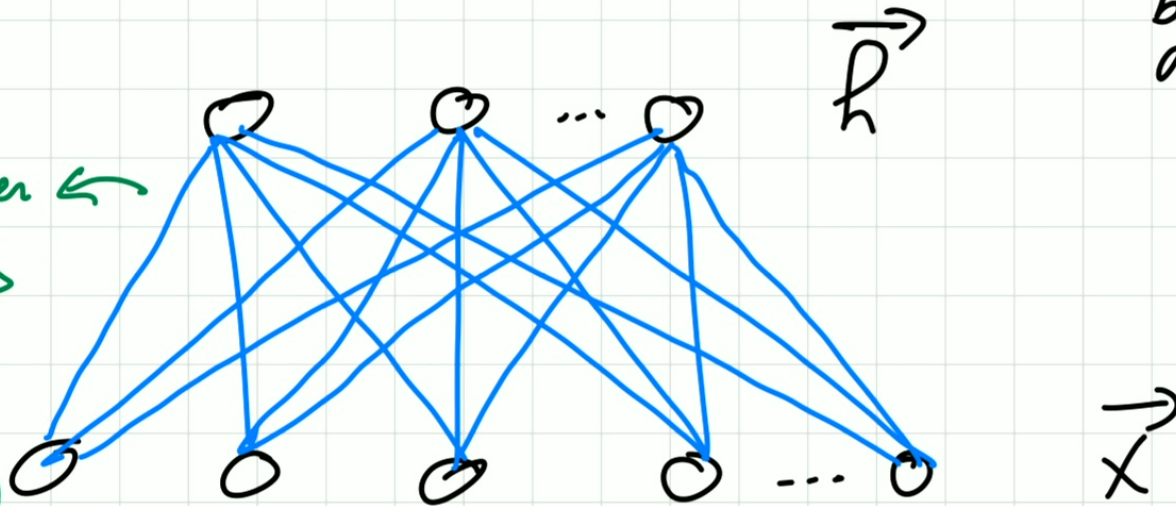
→ BM is systematically improvable with the increase of the number of hidden variables $\leftarrow \begin{matrix} \rightarrow \\ h \end{matrix} \right\rangle$

number of hidden variables " h "

.....

↳ Restricted Boltzmann Machine (RBM) → Also pioneered by Hinton and collaborators

Remove
intra-layer
connections
↓
For
Easier
analytical
computations



The hidden variables h_1, h_2, h_3 .

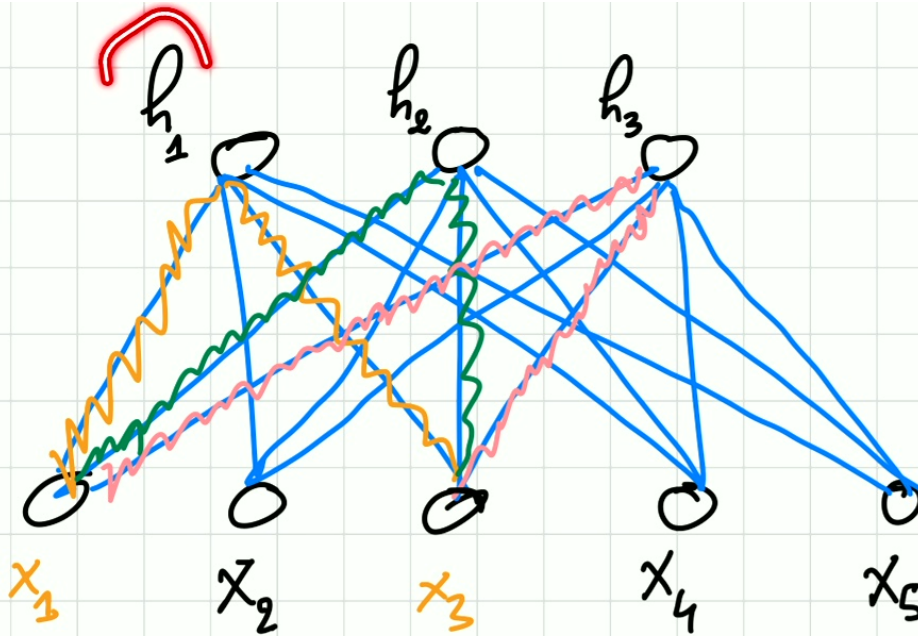
→ In more details:

~~$\tilde{W}_{ij} x_i x_j$~~ ~~$W'_{ij} h_i h_j$~~

$$E_{\vec{x}}(\vec{x}, \vec{h}) = - \sum_{ij} W_{ij} x_i h_j - \sum_{i=1}^N b_i x_i - \sum_{j=2}^M c_j h_j \dots$$

and $M = \lfloor \alpha N \rfloor$ (Number of hidden variables)
 Hyperparameter $\in \mathbb{R}^+$

Main idea:
≈ ≈



For instance:

→ Correlation between x_1 and x_3 is introduced through the hidden variables h_1, h_2, h_3 .

the hidden variables h_1, h_2, h_3 .

→ In more details:

~~$\tilde{W}_{ij} x_i h_j$~~ ~~$W_{ij} h_i h_j$~~

$$E_{\vec{x}}(\vec{x}, \vec{h}) = - \sum_{ij} W_{ij} x_i h_j - \sum_{i=1}^N b_i x_i - \sum_{j=1}^M c_j h_j$$

and $M = \lfloor \alpha N \rfloor$ (Number of hidden variables)
Hyperparameter $\in \mathbb{R}^+$

$$P(\vec{x}) = \sum P(\vec{x}, \vec{h})$$

$$E_{\vec{x}}(\vec{x}, \vec{h}) = - \sum_{ij} W_{ij} x_i h_j - \sum_{i=1}^N b_i x_i - \sum_{j=1}^M c_j h_j$$

and $M = \lfloor \alpha N \rfloor$ (Number of hidden variables)

Hyperparameter $\in \mathbb{R}^+$

$$P_{\lambda}(\vec{x}) = \sum_{\vec{h}} P(\vec{x}, \vec{h})$$

Marginalization

$$= \frac{\exp(-E_{\lambda}(\vec{x}, \vec{h}))}{Z_{\lambda}}$$

the hidden variables h_1, h_2, h_3 .

→ In more details:

~~$\tilde{W}_{ij} x_i x_j$~~ ~~$W'_{ij} h_i h_j$~~

$$E_{\vec{x}}(\vec{x}, \vec{h}) = - \sum_{ij} W_{ij} x_i h_j - \sum_{i=1}^N b_i x_i - \sum_{j=2}^M c_j h_j$$

and $M = \lfloor \alpha N \rfloor$ (Number of hidden variables)
Hyperparameter $\in \mathbb{R}^+$

$$n(\rightarrow) \leftarrow n(\rightarrow) \vec{n}$$

* Industry applications

→ Recommendation engines (collaborative filtering ~ 2000)

→ Netflix competition (2006)

→ Unsupervised pre-training for deep learning (pre-2012)
Alex Net

* Many-body Physics applications:

*1st application: Train RBM on data drawn from $P(\vec{x})$
↓
unknown
so that: $P_{\text{RBM}}(\vec{x}) \approx P(\vec{x})$

↳ After training, we can sample RBM to
compute observables $\langle Q \rangle$

→ In more details:

$$\cancel{W_{ij} x_i x_j} \quad \cancel{W_{ij} h_i h_j}$$

$$E_{\vec{x}}(\vec{x}, \vec{h}) = - \sum_{ij} W_{ij} x_i h_j - \sum_{i=1}^N b_i x_i - \sum_{j=1}^M c_j h_j$$

and $M = \lfloor \alpha N \rfloor$ (Number of hidden variables)
Hyperparameter $\in \mathbb{R}^+$

$$P_{\lambda}(\vec{x}) = \sum_{\vec{h}} P(\vec{x}, \vec{h})$$

Marginalization

$$= \frac{\exp(-E_{\lambda}(\vec{x}, \vec{h}))}{Z_{\lambda}}$$

*1st application: Train RBM on data drawn from $P(\vec{x})$
↓
Unknown

so that: $P_{\text{RBM}}(\vec{x}) \approx P(\vec{x})$

↳ After training, we can sample RBM to

compute observables $\langle Q \rangle$

$$\langle Q \rangle = \sum_{\vec{x}} P_{\text{RBM}}(\vec{x}) Q(\vec{x})$$

$$\langle Q \rangle = \frac{1}{Z} \sum_{\vec{x}} \sum_{\vec{h}} \exp(-E_{\vec{x}, \vec{h}}) Q(\vec{x})$$

Kullback-Liebler Divergence

$$D_{KL}(P \parallel P_{RBM}) = + \sum_{\vec{x}} P(\vec{x}) \log \left(\frac{P(\vec{x})}{P_{RBM}(\vec{x})} \right)$$

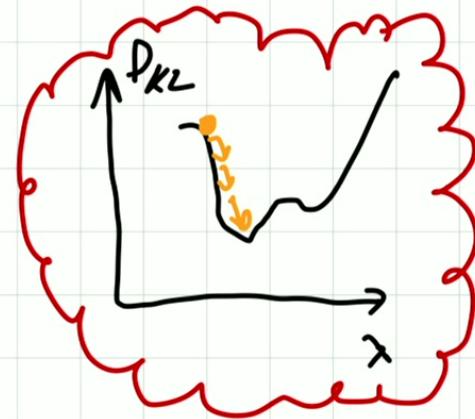
Cost function

Property:

$$D_{KL} = 0 \text{ iff } P = P_{RBM} \text{ and } D_{KL} > 0 \text{ (otherwise)}$$

↳ We can perform gradient descent:

$$\vec{\lambda} \leftarrow \vec{\lambda} - \eta \partial_{\vec{\lambda}} D_{KL}$$



$$D_{KL}(P \parallel P_{RBM}) = - \sum_{\vec{x}} P(\vec{x}) \log \left(\frac{P(\vec{x})}{P_{RBM}(\vec{x})} \right)$$

Cost
function

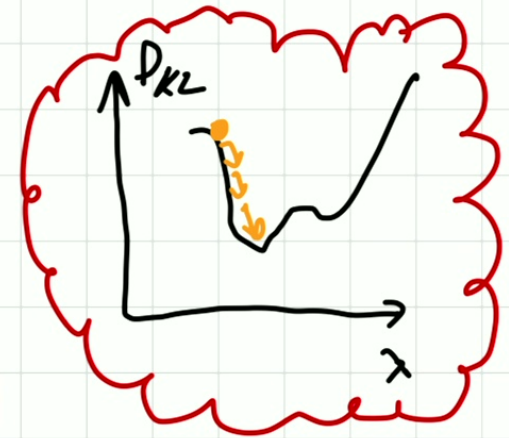
Property:

$D_{KL} = 0$ iff $P = P_{RBM}$ and $D_{KL} > 0$ (otherwise)

We can perform gradient descent:

Adam
optimizer

$$\vec{\lambda} \leftarrow \vec{\lambda} - \eta \partial_{\vec{\lambda}} D_{KL}$$



* Question: Is D_{KL} tractable?

... Hill ... of ...

$$D_{KL}(P \parallel P_{RBM}) = + \sum_{\vec{x}} P(\vec{x}) \log \left(\frac{P(\vec{x})}{P_{RBM}(\vec{x})} \right)$$

Cost
function

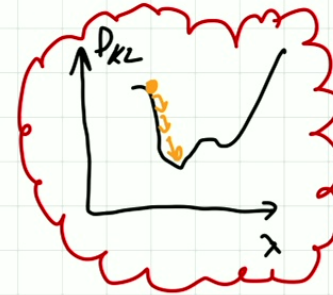
Property:

$$D_{KL} = 0 \text{ iff } P = P_{RBM} \text{ and } D_{KL} > 0 \text{ (otherwise)}$$

↳ We can perform gradient descent:

"Adam"
optimizer

$$\eta \leftarrow \eta - \eta \partial_{\eta} D_{KL}$$



* Question: Is D_{KL} tractable?

$$D_{KL}(P \parallel P_{RBM}) = \underbrace{\sum_{\vec{x}} P(\vec{x}) \log(P(\vec{x}))}_{-H_{data} = -\text{entropy of data}} - \sum_{\vec{x}} P(\vec{x}) \log(P_{RBM}(\vec{x}))$$

Cert
function

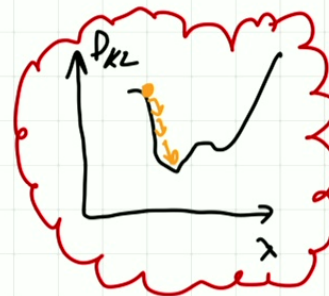
→ convex

$$D_{KL} = 0 \text{ iff } P = P_{RBM} \text{ and } D_{KL} > 0 \text{ (otherwise)}$$

↳ We can perform gradient descent:

↳ "Adam"
optimizer

$$\lambda \leftarrow \lambda - \eta \partial_{\lambda} D_{KL}$$



* Question: Is D_{KL} tractable?

$$D_{KL}(P \parallel P_{RBM}) = \sum_{\vec{x}} P(\vec{x}) \log(P(\vec{x}))$$

$-H_{data} = -\text{entropy of data}$

$$- \sum_{\vec{x}} P(\vec{x}) \log(P_{RBM}(\vec{x}))$$

Independent
of parameters λ

$$\approx -H_{data} - \frac{1}{|D|} \sum_{\vec{x} \in D} \log(P_{RBM}(\vec{x}))$$

→ Recommendation engines (Collaborative filtering ~ 2000)

→ Netflix competition (2006)

→ Unsupervised pre-training for deep learning (Pre-2012)
AlexNet

* Many-body Physics applications:

*1st application: Train RBM on data drawn from $P(\vec{x})$
↓
Unknown
so that: $P_{\text{RBM}}(\vec{x}) \approx P(\vec{x})$

↳ After training, we can sample RBM to
compute observables $\langle Q \rangle$

$$\langle Q \rangle = \sum_{\vec{x}} P_{\text{RBM}}(\vec{x}) Q(\vec{x})$$

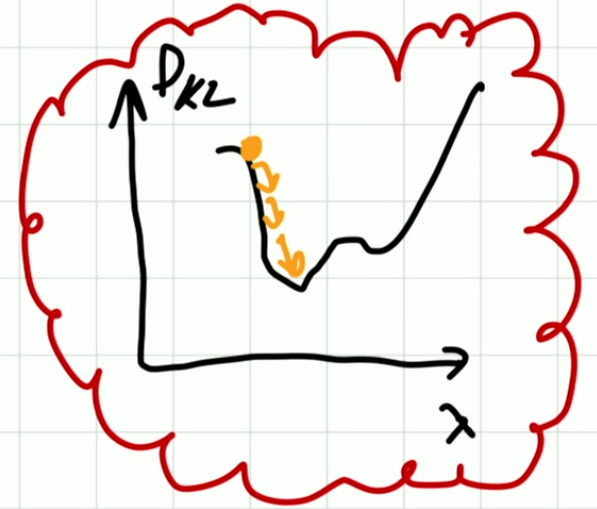
$$\langle Q \rangle = \frac{1}{Z} \sum_{\vec{x}} \sum_{\vec{h}} \exp(-E_{\vec{x}, \vec{h}}) Q(\vec{x})$$

Condition \rightarrow Property: $D_{KL} = 0$ iff $P = P_{RBM}$ and $D_{KL} > 0$ (otherwise)

We can perform gradient descent:

Adam optimizer

$$\lambda \leftarrow \lambda - \eta \partial_{\lambda} D_{KL}$$



Question: Is D_{KL} tractable?

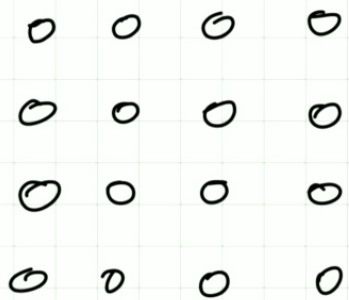
-HIL - - extension of

$$+ \frac{1}{|D|} \sum_{\vec{x} \in D} z_{\vec{x}} z_{\vec{x}'} \quad \text{Intractable} \rightarrow \text{Contradictory Divergence (CD)}$$

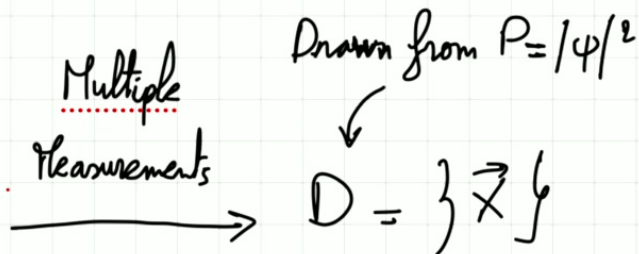
Negative phase (See additional Notes)

Demo:
 → See tutorial tomorrow (December for Rydberg atoms)

Example:



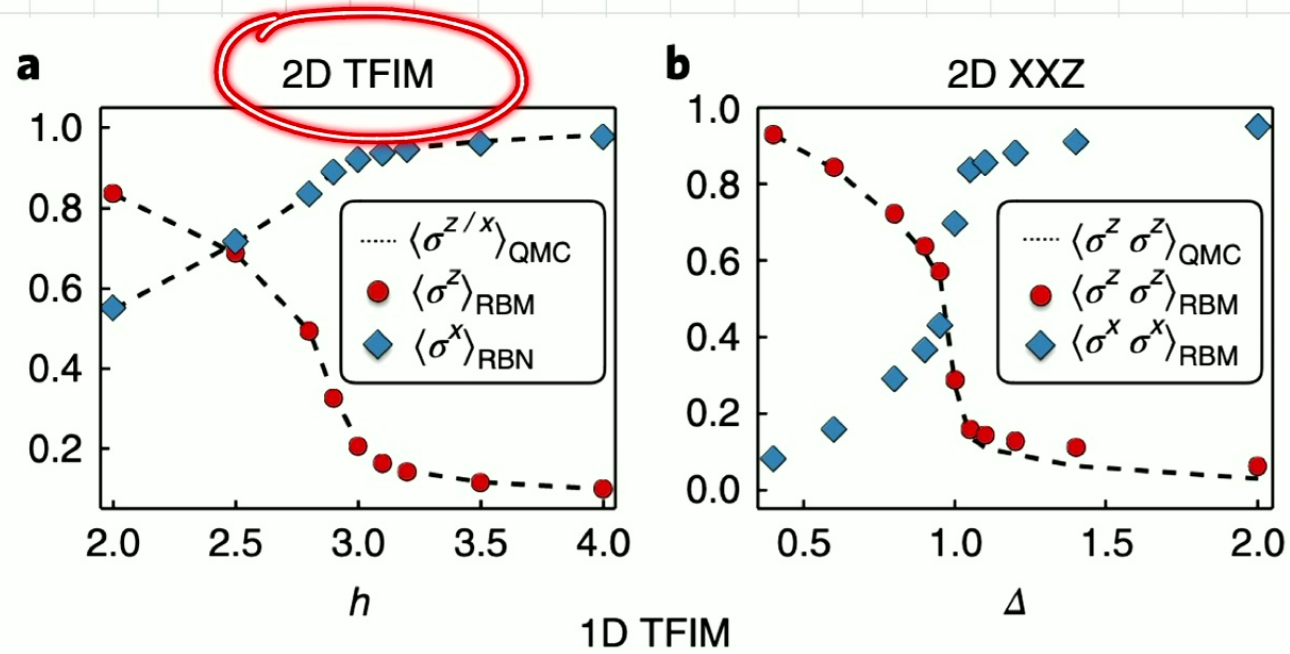
Quantum system $|\psi\rangle$
 (e.g. Rydberg atoms)



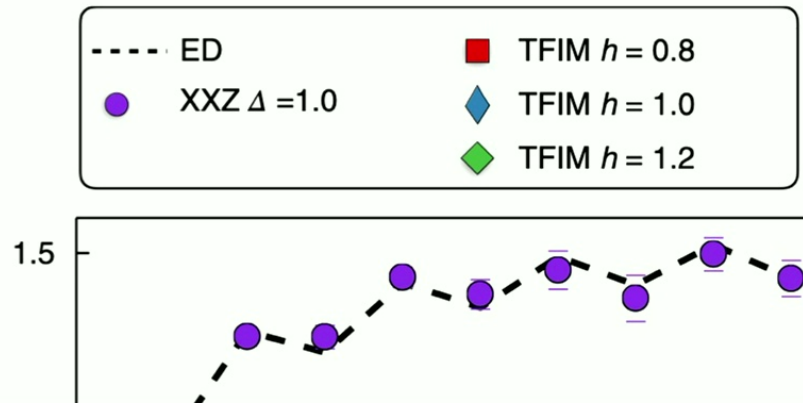
↓
 Train an RBM

$$\Psi_{\text{RBM}} \simeq \Psi = \sqrt{p} e^{i\phi}$$

→ → → → →



See: "Neural Network
quantum state
 tomography",



* 2nd application: (Variational Monte Carlo - VMC)

→ Given a quantum system described by Hamiltonian

$$\hat{H} \text{ (e.g. } \hat{H} = - \sum_i \hat{\sigma}_i^z \hat{\sigma}_{i+1}^z - B_x \sum_i \hat{\sigma}_i^x \text{)}$$

* We want to approximate the ground state: $|\psi_G\rangle = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{pmatrix} \in \mathbb{C}^{2^N}$

$$\hat{H} |\psi_G\rangle = E_G |\psi_G\rangle$$

↳ We can use an RBM with

Tractable
but not
normalized

$$\psi_{\text{RBM}}(\vec{x}) = \sum_{\vec{h}} e^{-E_{\vec{x}}(\vec{x}, \vec{h})}$$

one could use
complex-valued
parameters

... $|\hat{H}|(\vec{x}) \rangle$

↳ We can use an RBM with

Tractable
but not
normalized

$$\Psi_{\text{RBM}}(\vec{x}) = \sum_{\vec{h}} e^{-E_{\vec{x}}(\vec{x}, \vec{h})}$$

one could use
complex-valued
parameters

↳ Train RBM using $\frac{\langle \Psi_{\text{RBM}} | \hat{H} | \Psi_{\text{RBM}} \rangle}{\langle \Psi_{\text{RBM}} | \Psi_{\text{RBM}} \rangle}$

as a cost function.

↳ See tutorial with NetKet

<https://netket.readthedocs.io/en/v3.12.0/tutorials/gs-heisenberg.html>

↳ See Carleo & Troyer, Science 2017

↳ Open quantum systems: Vicentini et al.
arXiv: 1902.10104 2019

↳ "Restricted Boltzmann Machines in Quantum Physics", Nature
2019

* Main limitations: (Motivating RNNs & Transformers)

* Block-Gibbs sampling is not perfect.

$$\begin{aligned} * \underbrace{\partial_{\vec{\lambda}} D_{KL}}_{\substack{\text{Gradients estimation} \\ \text{is challenging}}} &\approx -\frac{1}{|D|} \sum_{\vec{x} \in D} \log \left(\overbrace{\sum_{\vec{h}} \exp(-E_{\vec{\lambda}}(\vec{x}, \vec{h}))}^{\text{Tractable}} \right) \\ &+ \underbrace{\frac{1}{|D|} \sum_{\vec{x} \in D} \partial_{\vec{\lambda}} z_{\vec{x}}}_{\substack{\text{Intractable} \rightarrow \text{Contradictory} \\ \text{Divergence (CD)}}} \end{aligned}$$

⚠