

Restricted Boltzmann Machines

Patrick Huembeli

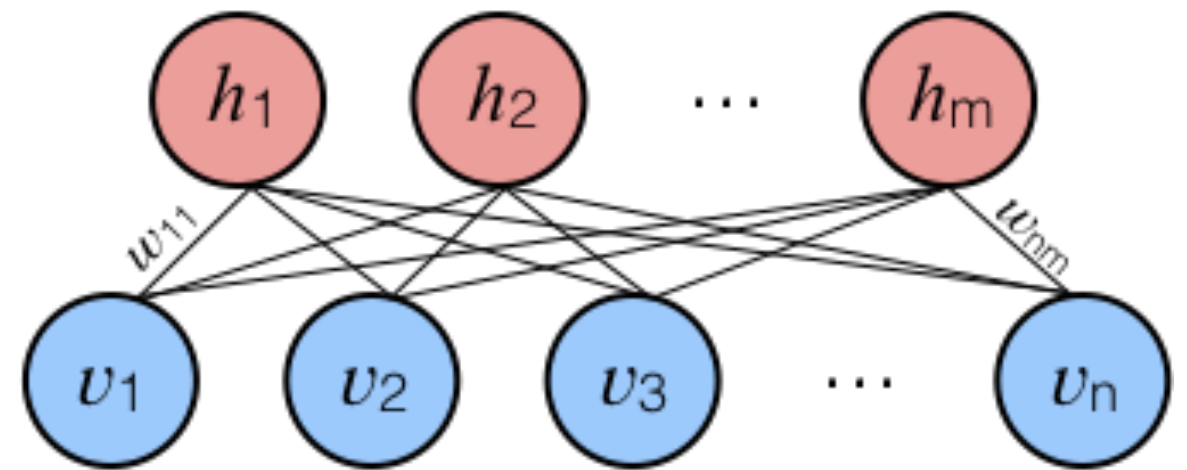
Build an RBM

Ingredients:

- Bipartite graph: hidden and visible
- Binary units $\{0,1\}$ **or** $\{-1,1\}$
- Probabilistic model
 - Normal Neural Networks are deterministic
 - We want to find e.g.

$$p(v_i = 1)$$

$$p(h_i = 0)$$



<https://medium.com/@MeTroFuN/python-mxnet-tutorial-1-restricted-boltzmann-machines-using-ndarray-f77578648ecf>

Parametrise weights and biases with:

$$\lambda = (W, b, c)$$

Build an RBM

$$E_{\lambda}(\mathbf{v}, \mathbf{h}) = - \sum_{i,j} w_{ij} v_i h_j - \sum_i b_i v_i - \sum_j c_j h_j$$

Boltzmann distribution

$$p_{\lambda}(\mathbf{v}, \mathbf{h}) = e^{-E_{\lambda}(\mathbf{v}, \mathbf{h})} / Z$$

Partition Function

$$Z = \sum_{\mathbf{v}, \mathbf{h}} p_{\lambda}(\mathbf{v}, \mathbf{h})$$

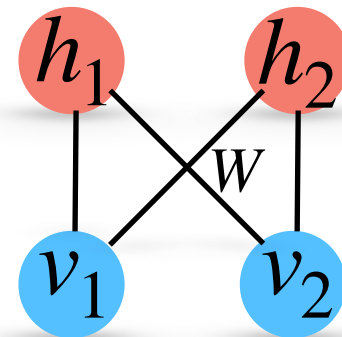
$$\rightarrow p_{\lambda}(\mathbf{v}) = \sum_{\mathbf{h}} p_{\lambda}(\mathbf{v}, \mathbf{h}) = e^{-\mathcal{E}(\mathbf{v})} / Z$$

$$\mathcal{E}(\mathbf{v}) = - \sum_i b_i v_i - \sum_j \text{softplus}(c_j + \sum_i w_{i,j} v_i)$$

$$\text{softplus}(x) = \ln(1 + \exp(x))$$

Min. Energy gives max. Probability !!!

Example:



$$b_i = c_i = 0$$

$$w_{i,j} = 1, \forall i, j$$

$$E = -v_1 h_1 - v_1 h_2 - v_2 h_1 - v_2 h_2$$

Which configuration is most likely??

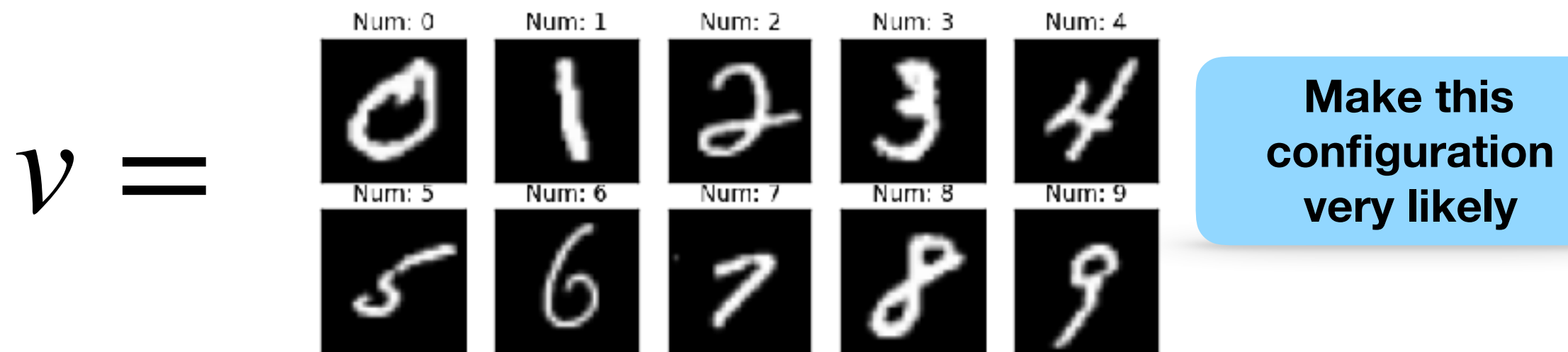
Homework:

**Calculate Probabilities for this case for all configurations.
To do so you will have to calculate the partition function
and the energies for all possible configurations.**

Do the same with $w_{i,j} = -1, \forall i, j$

Why is this useful?

- RBM can learn probability distributions $p_{\lambda}(\mathbf{v})$ by finding the parameters that:



Make this configuration very unlikely

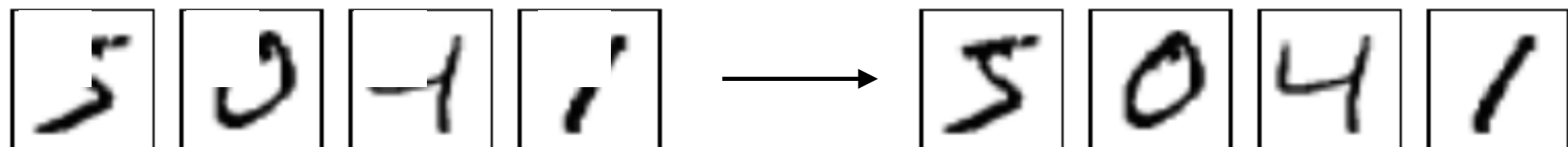


Make this configuration quite likely

- Learning is unsupervised, no labels needed

Why is this useful?

- We can sample from the learned distribution via Gibbs sampling
- Get new configurations
- Generate new data
- Restore 'defect' data (Homework)



1st Summary

- RBMs are probabilistic models
- We here only look at binary RBMs, which means all nodes are either 0 or 1
- We define energy and adjust parameters such that the configurations of our data have the smallest energy
 - Which is the highest probability
 - This is inspired by statistical physics (e.g. classical Spins)
 - Preferred configurations have lowest energy
- Think of training data as probability distribution.
 - Distribution of pixels that we want to learn

Build an RBM

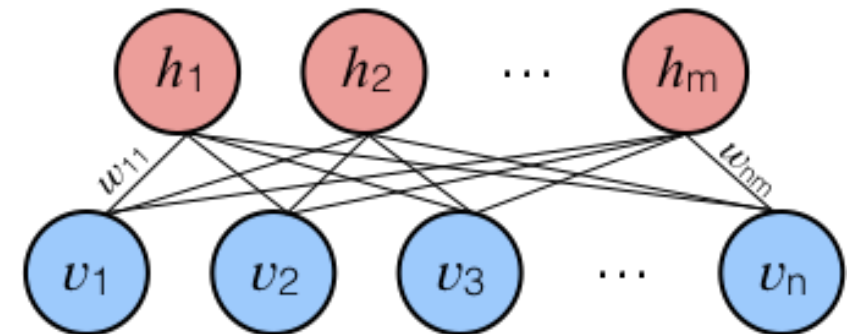
Probabilities:

$$p_{\lambda}(\mathbf{v}, \mathbf{h}) = p_{\lambda}(\mathbf{v} | \mathbf{h}) p_{\lambda}(\mathbf{h})$$

$$p_{\lambda}(\mathbf{h}, \mathbf{v}) = p_{\lambda}(\mathbf{h} | \mathbf{v}) p_{\lambda}(\mathbf{v})$$

Because there are no connections between the units of the same layer, the conditional distributions factorise:

$$p_{\lambda}(\mathbf{v} | \mathbf{h}) = \prod_i p_{\lambda}(v_i | \mathbf{h})$$
$$p_{\lambda}(\mathbf{h} | \mathbf{v}) = \prod_j p_{\lambda}(h_j | \mathbf{v})$$



Build an RBM

- Calculate conditional probabilities

$$p_{\lambda}(\mathbf{v}_i = 1 \mid \mathbf{h}) = \mathcal{S} \left(b_i + \sum_j h_j w_{i,j} \right)$$

$$p_{\lambda}(\mathbf{h}_i = 1 \mid \mathbf{v}) = \mathcal{S} \left(c_j + \sum_i v_j w_{i,j} \right)$$

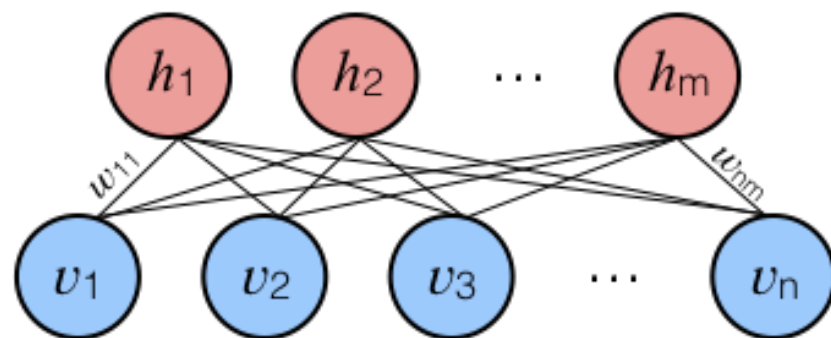
$$\mathcal{S}(x) = \frac{1}{1 + e^{-x}}$$

Derivation:

<http://www.iro.umontreal.ca/~bengioy/ift6266/H14/ftml-sec5.pdf>

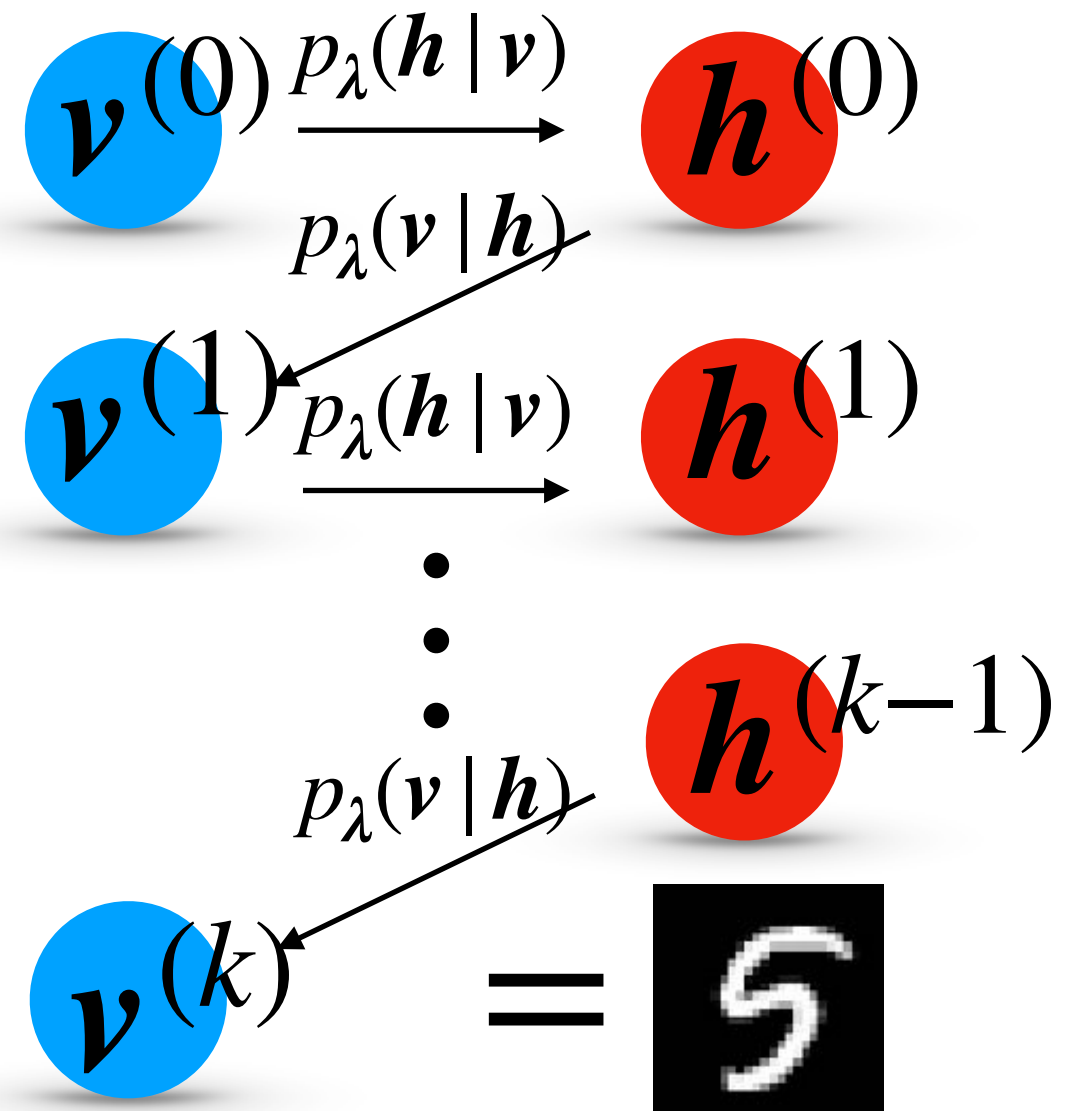
Build an RBM: Example

- After RBM is trained on MNIST, the configurations of MNIST are most likely!



Input v is given

$v^{(0)}$



Markov chain converges to stationary solution!
Which means low energy solution

Train an RBM

‘Cost function’:

$$C_{\lambda} = D_{KL}(q || p_{\lambda}) = \sum_{\mathbf{v}} q(\mathbf{v}) \log \left(\frac{q(\mathbf{v})}{p_{\lambda}(\mathbf{v})} \right)$$

q

is the data distribution, for example MNIST. How are the pixels normally arranged in a MNIST data set? What is the probability $q(\mathbf{v})$ of a configuration \mathbf{v} in your data set?

p_{λ}

is the distribution learned by the RBM. The parameters λ determine the probability of a certain configuration \mathbf{v} .

Train an RBM

- Minimize Kullback-Leibler (KL) divergence
- To calculate gradient of KL-Divergence we need to calculate expectation values over the model distribution.

Problem!

$$\nabla_{\lambda} C_{\lambda} \approx \langle \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) \rangle_{\mathcal{D}} - \langle \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) \rangle_{p_{\lambda}}$$

$$p_{\lambda}(\mathbf{v}) = e^{-\mathcal{E}(\mathbf{v})}/Z$$

- Partition function Z generally not accessible
- Derivation from : https://qucumber.readthedocs.io/en/stable/static/RBM_tutorial.pdf

Train RBM

- Derivation on https://github.com/PatrickHuembeli/QML-Course-UPC-2018/blob/master/RBM_Slides_and_Notes/RBM_gradient_derivation.pdf

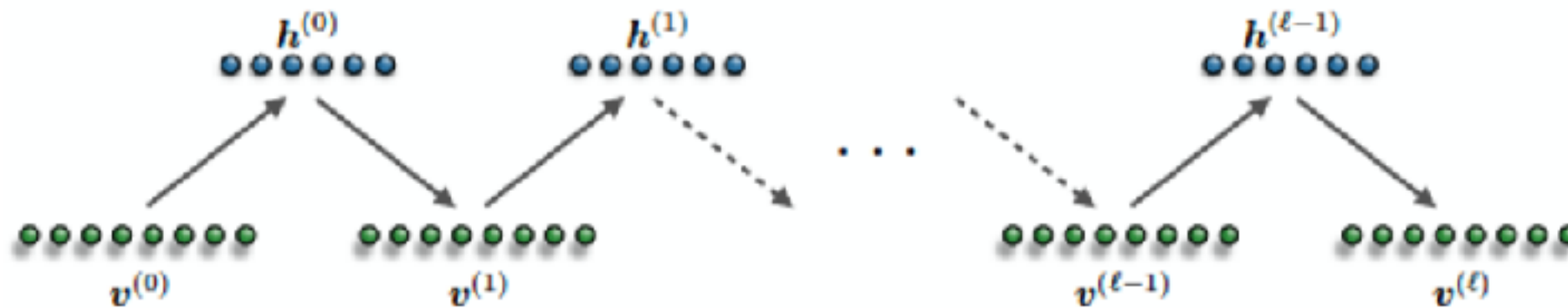
Train RBM

- The problem is we cannot calculate $p_{\lambda}(\mathbf{v})$
- Because to calculate $Z = \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h})$ sum grows exponentially
- Approximate expectation value with estimator sampled from model using Gibbs sampling

$$\langle \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) \rangle_{p_{\lambda}} = \sum_{\mathbf{v}} p_{\lambda}(\mathbf{v}) \mathcal{E}_{\lambda}(\mathbf{v}) \approx \frac{1}{M} \sum_{\mathbf{v}^{(l)} \in \mathcal{M}} \mathcal{E}_{\lambda}(\mathbf{v}^{(l)})$$

Gibbs sampling?

- Ideally Markov chain converges to stationary solution



- If we sample back and forth, we approximate the distribution of the RBM
- Initial vector is a training sample (Homework)
- Empirics show, that one Gibbs step is enough

Contrastive divergence

- Instead of taking the model distribution $p_\lambda(\mathbf{v})$
- We calculate the expectation value with an estimator

$$\langle \nabla_\lambda \mathcal{E}_\lambda(\mathbf{v}) \rangle_{p_\lambda} = \sum_{\mathbf{v}} p_\lambda(\mathbf{v}) \mathcal{E}_\lambda(\mathbf{v}) \approx \frac{1}{M} \sum_{\mathbf{v}^{(k)} \in \mathcal{M}} \mathcal{E}_\lambda(\mathbf{v}^{(k)})$$

- The number of Gibbs steps we take is k
- The more steps we take the closer we converge to the actual model distribution
- But it also takes more time
- Contrastive divergence is often referred to as CD-k, where k gives the number of steps

Calculate gradients explicitly

$$\frac{\partial \mathcal{E}_\lambda(\mathbf{v})}{\partial w_{i,j}} = -v_i \cdot p_\lambda(h_j = 1 | \mathbf{v})$$

$$\frac{\partial \mathcal{E}_\lambda(\mathbf{v})}{\partial c_j} = -p_\lambda(h_j = 1 | \mathbf{v})$$

$$\frac{\partial \mathcal{E}_\lambda(\mathbf{v})}{\partial b_i} = -v_i$$

Check as Homework!

Parameter update

- Like in stochastic gradient descent $\lambda \leftarrow \lambda - \eta \nabla_{\lambda} C_{\lambda}$
- Reminder: $\nabla_{\lambda} C_{\lambda} \approx = \langle \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) \rangle_{\mathcal{D}} - \langle \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) \rangle_{p_{\lambda}}$
- E.g. weight update

$$w_{i,j} \leftarrow w_{i,j} - \eta \left\langle \frac{\partial \mathcal{E}_{\lambda}(\mathbf{v})}{\partial w_{i,j}} \right\rangle_{\mathcal{D}} + \eta \left\langle \frac{\partial \mathcal{E}_{\lambda}(\mathbf{v})}{\partial w_{i,j}} \right\rangle_{p_{\lambda}}$$

$$w_{i,j} \leftarrow w_{i,j} + \eta \langle v_i \cdot p_{\lambda}(h_j = 1 | \mathbf{v}) \rangle_{\mathcal{D}} - \eta \langle v_i \cdot p_{\lambda}(h_j = 1 | \mathbf{v}) \rangle_{p_{\lambda}}$$

$$w_{i,j} \leftarrow w_{i,j} + \eta \frac{1}{|\mathcal{D}|} \sum_{\mathbf{v}_i \in \mathcal{D}} v_i \cdot p_{\lambda}(h_j = 1 | \mathbf{v}) - \eta \frac{1}{|\mathcal{M}|} \sum_{\mathbf{v}_i^{(k)} \in \mathcal{M}} v_i \cdot p_{\lambda}(h_j = 1 | \mathbf{v})$$

**Sum over batch
from data**

**Sum over batch of
Gibbs samples**

Restricted Boltzmann Machines in Physics

RBM as ansatz for wavefunctions

$$\Psi = \sum_{\sigma} \Psi(\sigma) |\sigma\rangle \quad P(\sigma) \propto |\Psi(\sigma)|^2 \quad \psi_{\lambda}(\sigma) = \sqrt{\frac{p_{\lambda}(\sigma)}{Z_{\lambda}}}$$
$$\Psi(\sigma) = \langle \sigma | \Psi \rangle$$

- Variational Monte Carlo and minimize e.g. energy (Troyer et al.)
- Learn state from samples / measurements
- Ansatz is efficient.
 - ➡ Grows polynomial in system size

RBM as ansatz for wavefunctions

Wavefunction of physical system

$$\Psi(\sigma) = \langle \sigma | \Psi \rangle$$

Take RBM ansatz

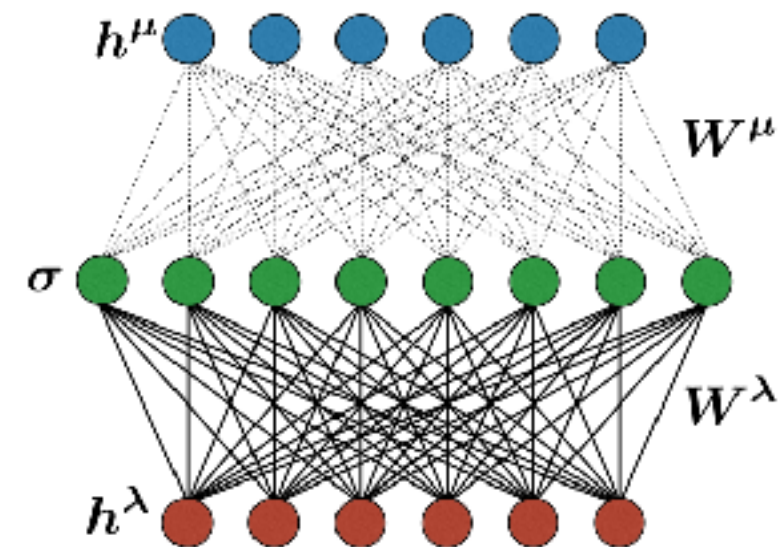
Real-positive wavefunction

Complex wavefunction

$$\psi_{\lambda}(\sigma) = \sqrt{\frac{p_{\lambda}(\sigma)}{Z_{\lambda}}}$$

$$\psi_{\lambda,\mu}(\sigma) = \sqrt{\frac{p_{\lambda}(\sigma)}{Z_{\lambda}}} e^{i\phi_{\mu}(\sigma)}$$

$$\phi_{\mu}(\sigma) = \log p_{\mu}(\sigma)$$



State tomography

- For phases we need to change basis.
- Make measurements on system in different basis
 $b = 0, 1, \dots$

$$P_b(\boldsymbol{\sigma}^{[b]}) \propto \left| \Psi(\boldsymbol{\sigma}^{[b]}) \right|^2$$

- Learn this distribution with RBM

$$\psi_{\lambda,\mu}(\boldsymbol{\sigma}^{[b]}) = \sum_{\boldsymbol{\sigma}} U_b(\boldsymbol{\sigma}, \boldsymbol{\sigma}^{[b]}) \psi_{\lambda,\mu}(\boldsymbol{\sigma})$$

- In contrast to positive-real wavefunction, the unitaries will be part of the training.

Results

- RBMs are efficient representation for quantum states
- We can sample from it
- We can calculate any expectation value via estimators
- Open source software for QST

➡ <https://github.com/PIQuIL/QuCumber>

