

Our previous forays into unsupervised learning have been concerned with finding structure in unlabelled data (e.g. dimensional reduction, t-SNE, clustering).

Now we discuss another important flavor, generative modelling, which is concerned with learning the underlying probability distribution that describes some observed data set.

### Generative Modelling

Consider the data set of observations  $\tilde{X} = \{\tilde{x}_n\}_{n=1}^N$ , where  $\tilde{x}_n \in \mathbb{R}^D, \mathbb{Z}_2^D, \dots$  (could also be complex if our data set consists of coefficients of quantum wavefunctions). There exists an underlying probability distribution  $P(\tilde{x})$  that describes the data (i.e. it is drawn from it) that we want to discover.

Goal: Introduce a parameterization (model) of the distribution  $P_\lambda(\tilde{x})$  that depends on some parameters  $\vec{\lambda}$  to find the parameters such that  $P_\lambda(\tilde{x}) \approx P(\tilde{x})$ .

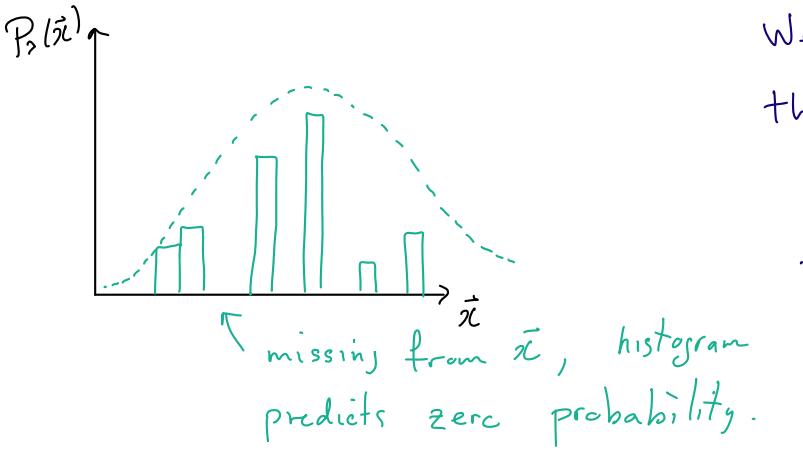
$\underbrace{P_\lambda(\tilde{x})}_{\substack{\text{model} \\ \text{dist.}}}$ 

 $\underbrace{P(\tilde{x})}_{\substack{\text{data} \\ \text{dist.}}}$

Naively with enough observations:

$$P_\lambda(\tilde{x}) = P(\tilde{x}) = \frac{1}{N} \sum_{n=1}^N \delta_{\tilde{x}_n, \tilde{x}}$$

However unless  $N$  is much larger than the total number of possible observations, this will fail to generalize.



We want to learn a  $P_\theta(\vec{x})$  that:

- i) Generalizes well
- ii) Can be sampled to generate observations not seen in the data set.

Idea: Use a type of parameterization, energy based model that can be trained to minimize the Kullback-Liebler (KL) divergence

### Energy Based Models

Let us restrict our set of observations to binary data, i.e.  $\vec{x}_n \in \mathbb{Z}_2^D$ . E-B models can then be thought of in terms of Ising Models that are a graphical representation of the probability distribution.

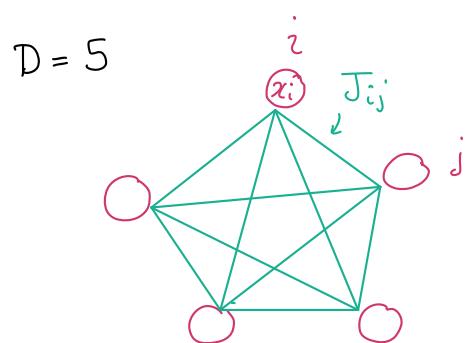
Quick history:

#### 1982 Hopfield Network

$$E(\vec{x}) = - \sum_i a_i x_i - \frac{1}{2} \sum_{i,j} x_i J_{ij} x_j \quad ; \quad x_i = \{0,1\}$$

↑                      ↓  
configuration      magnetic field      spins      couplings

Here our parameters are  $\lambda = \{\bar{a}, J\}$



We take the probability as a Boltzmann Distribution

$$P_p(\vec{x}) = \frac{1}{Z} e^{-E(\vec{x})} \quad \text{with} \quad Z = \sum_{\vec{x}} e^{-E(\vec{x})} \quad (k_B T = 1)$$

↑  
sum over all  $2^D$  configurations

$$\begin{aligned} \text{i.e. } \sum_{\vec{x}} &\equiv \sum_{x_1=0}^1 \cdot \sum_{x_2=0}^1 \cdots \sum_{x_D=0}^1 \\ &= \prod_{i=1}^D \sum_{x_i=0}^1 \end{aligned}$$

Unfortunately the "all-to-all" coupling  $J_{ij}$  makes these models impossible to train. However there is a trick:

ii) write  $J_{ij} = \sum_k w_{ik} w_{kj}$  (always possible via SVD)

iii) introduce auxiliary (latent) variables  $h_j$   $\downarrow$  perform a Hubbard-Stratonovich transformation, assuming they are continuous  $\downarrow$  normally distributed, this yields

$$P_p(\vec{x}) = \frac{\int d\vec{h} e^{-E(\vec{x}, \vec{h})}}{Z} \quad \begin{matrix} \text{No connections} \\ \text{within a layer!} \end{matrix}$$

$$E(\vec{x}, \vec{h}) = - \sum_i a_i x_i + \frac{1}{2} \sum_j h_j^2 - \sum_{i,j} x_i w_{ij} h_j$$

$\downarrow$  define a new energy based model

$$P_p(\vec{x}, \vec{h}) \equiv \frac{e^{-E(\vec{x}, \vec{h})}}{Z'} \quad ; \quad Z' = \sum_{\vec{x}, \vec{h}} e^{-E(\vec{x}, \vec{h})}$$

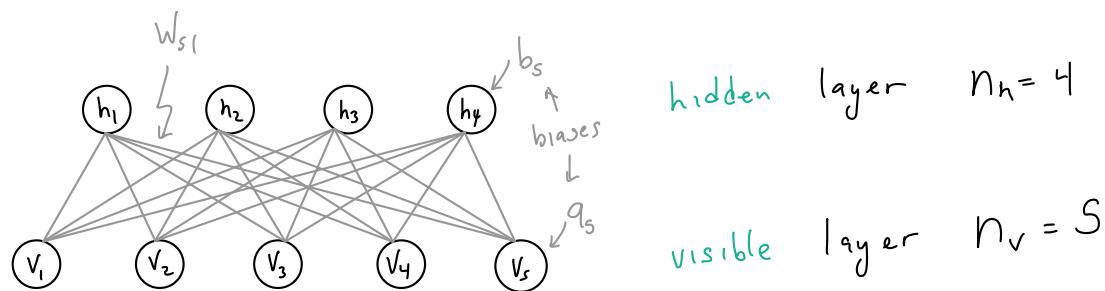
## Restricted Boltzmann Machines (RBMs)

Introduced by Hinton in 1986, these represent a trainable class of energy-based models with an energy:

$$E(\vec{v}, \vec{h}) = - \sum_{i=1}^{n_v} a_i v_i - \sum_{j=1}^{n_h} b_j h_j - \sum_{i,j} v_i w_{ij} h_j$$

visible units                                      hidden units                              symmetric weight matrix

$v_i \in \{0,1\}$   
 $h_j \in \{0,1\}$



\* No connections inside a layer! This is what makes the model trainable!

As above the joint probability distribution has a Boltzmann form:

$$P_2(\vec{v}, \vec{h}) = \frac{e^{-E(\vec{h}, \vec{v})}}{Z_2} \quad ; \quad Z_2 = \sum_{\vec{h}, \vec{v}} e^{-E(\vec{h}, \vec{v})}$$

The quantity we are interested in is the marginal distribution

$$P_2(\vec{v}) = \sum_{\vec{h}} P_2(\vec{v}, \vec{h})$$

by Bayes theorem this can be written in terms of a conditional distribution

$$P_z(\vec{v} | \vec{h}) = \frac{P_z(\vec{v}, \vec{h})}{P_z(\vec{h})} \quad \begin{matrix} \text{joint} \\ \text{conditional} \end{matrix}$$

similarly :  $P_z(\vec{h} | \vec{v}) = \frac{P_z(\vec{v}, \vec{h})}{P_z(\vec{v})}$

Before discussing how to train, let's dig deeper into the properties of our graphical model:

$$\begin{aligned} Z_z P_z(\vec{v}) &= \sum_{\vec{h}} e^{-E(\vec{h}, \vec{v})} \\ &= \sum_{\vec{h}} e^{\sum_i a_i v_i + \sum_j b_j h_j + \sum_{ij} v_i w_{ij} h_j} \\ &= e^{\sum_i a_i v_i \left( \prod_{j=1}^{n_h} \sum_{h_j=0}^1 \right)} e^{\sum_j h_j \left( b_j + \sum_i v_i w_{ij} \right)} \\ &= e^{\sum_i a_i v_i \left( \prod_{j=1}^{n_h} \sum_{h_j=0}^1 e^{h_j z_j} \right)} \quad \text{used } e^{\sum_j \alpha_j} = \prod_j e^{\alpha_j} \\ &= e^{\sum_i a_i v_i \frac{n_h}{\prod_{j=1}^{n_h}} (1 + e^{z_j})} \quad * \text{only possible} \\ &\quad \text{because of restricted form!} \end{aligned}$$

Now we can examine the conditional distribution:

$$\begin{aligned} P_z(\vec{h} | \vec{v}) &= \frac{P_z(\vec{v}, \vec{h})}{P_z(\vec{v})} = \frac{e^{-E(\vec{v}, \vec{h})}}{Z_z P_z(\vec{v})} = \frac{e^{\sum_i a_i v_i + \sum_j b_j h_j + \sum_{ij} v_i w_{ij} h_j}}{e^{\sum_i a_i v_i} \prod_j (1 + e^{z_j})} \\ &= \frac{\prod_j e^{(b_j + \sum_i v_i w_{ij}) h_j}}{\prod_j (1 + e^{z_j})} = \prod_j \frac{e^{z_j h_j}}{1 + e^{z_j}} \end{aligned}$$

↗ conditional distribution factorizes!

$$\Rightarrow P_z(\vec{h} | \vec{v}) = \prod_j P_z(h_j | \vec{v}) \Rightarrow h_j \text{ are iid!}$$

Now, recall there are only two possible values,  $h_j \in \{0, 1\}$

$$\therefore P_z(h_j=1 | \vec{v}) = \frac{e^{z_j}}{1+e^{z_j}} = \frac{1}{1+e^{-z_j}} = \sigma(z_j)$$

↙ our familiar Sigmoid!  
 ↘ normalized.

$$\therefore P_z(h_j=0 | \vec{v}) = \frac{1}{1+e^{z_j}} = 1 - \sigma(z_j)$$

$\therefore$  we now have a mechanism to sample the hidden units based on a fixed visible configuration, i.e. calculate  $P_z(h_j=1 | \vec{v})$  by comparing with a uniformly distributed  $r \in U_{[0,1]}$

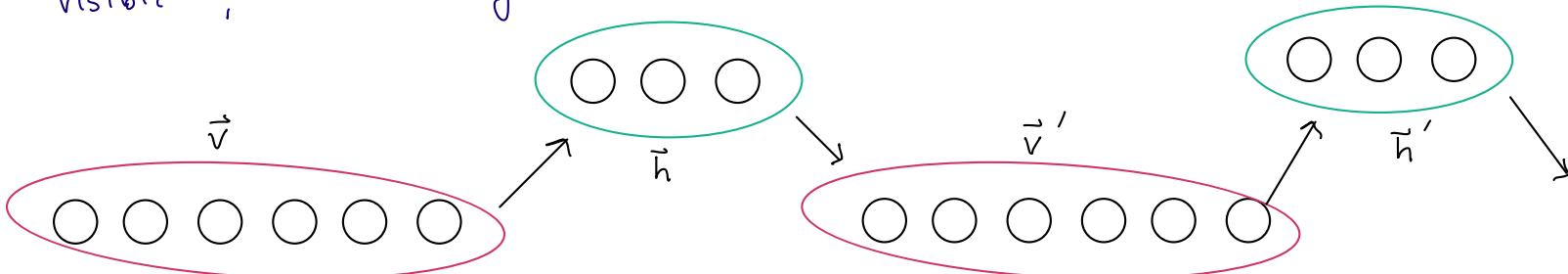
$$h_j = r < P_z.$$

Also, looking at the form of the energy, we can see that making the exchange  $a_i \rightarrow b_j$  is  $v_i \rightarrow h_j$  (since  $w_{ij} = w_{ji}$ )

$$P_z(v_i=1 | \vec{h}) = \sigma(z'_i) \quad \text{where} \quad z'_i = a_i + \sum_j w_{ij} h_j$$

$$P_z(v_i=0 | \vec{h}) = 1 - \sigma(z'_i)$$

This parameterization gives us the ability to sample both the visible & hidden layers with Markov Chain Monte Carlo:



After the MCMC has equilibrated, this gives us the ability to sample  $P_2(\vec{v})$  using Block Gibbs Sampling.

$$\left. \begin{aligned} P_2(\vec{h} \leftarrow \vec{v}) &= P_2(\vec{h} | \vec{v}) = \frac{P_2(\vec{v}, \vec{h})}{P_2(\vec{v})} \\ P_2(\vec{v} \leftarrow \vec{h}) &= P_2(\vec{v} | \vec{h}) = \frac{P_2(\vec{v}, \vec{h})}{P_2(\vec{h})} \end{aligned} \right\} \quad \frac{P_2(\vec{h} \leftarrow \vec{v})}{P_2(\vec{v} \leftarrow \vec{h})} = \frac{P_2(\vec{v})}{P_2(\vec{h})}$$

Satisfies detailed balance  
if stationary distributions  
are desired marginal ones!

### Training in the KL Divergence

Now we can return to the problem of choosing parameters  $\Delta = \{\vec{a}, \vec{b}, \vec{w}\}$  such that we can reproduce the data distribution  $P(\vec{v})$  using our training samples  $\{\vec{v}_n\}_{n=1}^N$ .  
 ↪  $\vec{v}$  is used instead of  $\vec{x}$ .  
 ↪ no  $\Delta$  here!

We want to minimize a cost function given by the Kullback-Liebler

divergence:

$$D_{KL}(P || P_2) = \sum_{\vec{v}} P(\vec{v}) \log \frac{P(\vec{v})}{P_2(\vec{v})}$$

Trace over all configs.

$$D_{KL} \geq 0$$

$$D_{KL} = 0 \Leftrightarrow P = P_2$$

$$= \underbrace{\sum_{\vec{v}} P(\vec{v}) \log P(\vec{v})}_{\text{- entropy of } P \equiv H} - \underbrace{\sum_{\vec{v}} P(\vec{v}) \log P_2(\vec{v})}_{\text{depends on parameters}} \quad \text{"binary cross-entropy"}$$

Unlike when training our DNN we can't just evaluate  $P(\vec{v})$  (this is unknown) but we can approximate the trace via

sampling over our training data:  $V = \{\vec{v}_n\}_{n=1}^N$

$$D_{KL}(P \parallel P_a) \simeq -H_V - \frac{1}{N} \sum_{\vec{v} \in V} \log P_a(\vec{v})$$

only this part depends  
on parameters  $\lambda$

*entropy of observations*

*log-likelihood from MLE  
for classification*

∴ Take our cost function:

$$\begin{aligned} C &= -\frac{1}{N} \sum_{\vec{v} \in V} \log P_a(\vec{v}) \\ &= -\langle \log P_a(\vec{v}) \rangle_{\text{data}} \end{aligned}$$

### Updating Parameters

With access to  $C$  we can update all model parameters  $\lambda = \{\vec{a}, \vec{b}, \vec{w}\}$  via algorithms like gradient descent. We need to take gradients wrt parameters:  $\lambda \leftarrow \lambda - \gamma \nabla_{\lambda} C$

Ultimately we will perform averages over mini-batches ; thus let us consider gradients evaluated at a fixed visible layer configuration  $\vec{v}$  (i.e. a data point).

Recall:  $P_a(\vec{v}) = \sum_{\vec{h}} P_a(\vec{v}, \vec{h})$

so

$$\begin{aligned} \frac{\partial \log P_a(\vec{v})}{\partial \lambda} &= \frac{\partial}{\partial \lambda} \left[ \log \sum_{\vec{h}} P_a(\vec{v}, \vec{h}) \right] \\ &= \frac{\partial}{\partial \lambda} \left[ \log \sum_{\vec{h}} \frac{e^{-E(\vec{v}, \vec{h})}}{Z_a} \right] \end{aligned}$$

$$\Rightarrow \frac{\partial}{\partial \lambda} \log P_{\lambda}(\vec{v}) = \frac{\partial}{\partial \lambda} \left[ \log \sum_{\vec{h}} e^{-E(\vec{v}, \vec{h})} \right] - \frac{\partial}{\partial \lambda} \left[ \log \sum_{\vec{h}, \vec{v}} e^{-E(\vec{v}, \vec{h})} \right]$$

$$= \frac{-1}{\sum_{\vec{h}} e^{-E(\vec{v}, \vec{h})}} \sum_{\vec{h}} e^{-E(\vec{v}, \vec{h})} \frac{\partial E}{\partial \lambda} + \frac{1}{\sum_{\vec{h}, \vec{v}} e^{-E(\vec{v}, \vec{h})}} \sum_{\vec{h}, \vec{v}} e^{-E(\vec{v}, \vec{h})} \frac{\partial E}{\partial \lambda}$$

Now recall our expression for the conditional probabilities:

$$P_{\lambda}(\vec{h} | \vec{v}) = \frac{P_{\lambda}(\vec{v}, \vec{h})}{P_{\lambda}(\vec{v})} = \frac{e^{-E(\vec{v}, \vec{h})}}{\cancel{\sum_{\vec{h}} e^{-E(\vec{v}, \vec{h})}}} = \frac{e^{-E(\vec{v}, \vec{h})}}{\sum_{\vec{h}} e^{-E(\vec{v}, \vec{h})}}$$

$$\Rightarrow \frac{\partial}{\partial \lambda} \log P_{\lambda}(\vec{v}) = - \sum_{\vec{h}} P_{\lambda}(\vec{h} | \vec{v}) \frac{\partial E}{\partial \lambda} + \sum_{\vec{h}, \vec{v}} P_{\lambda}(\vec{v}, \vec{h}) \frac{\partial E}{\partial \lambda}$$

so we can easily identify:

$$\frac{\partial E}{\partial a_i} = -v_i$$

$$\frac{\partial E}{\partial b_i} = -h_i$$

$$\frac{\partial E}{\partial w_{ij}} = -v_i h_j$$

as well as partially factorizing the last term using the probability product rule:

$$P_{\lambda}(\vec{v}, \vec{h}) = P_{\lambda}(\vec{h} | \vec{v}) P_{\lambda}(\vec{v})$$

So putting everything together:

$$\frac{\partial C}{\partial \lambda} = - \sum_{\vec{v}} P(\vec{v}) \frac{\partial}{\partial \lambda} \log P_{\lambda}(\vec{v})$$

$$\frac{\partial C}{\partial \alpha} = \sum_{\vec{v}} \left[ \sum_{\vec{h}} \underbrace{\frac{\partial E}{\partial \alpha} P_1(\vec{h} | \vec{v})}_{\text{Easy! Draw a single training sample, generate } \{\vec{h}\} \text{ from } \{\vec{v}\} \text{ using Block Gibbs, average over all samples } \vec{v}.} \right] P(\vec{v}) - \sum_{\vec{v}} \left[ \sum_{\vec{h}} \underbrace{\frac{\partial E}{\partial \alpha} P_2(\vec{h} | \vec{v})}_{\text{Hard! Need to average over the correct stationary distribution } P_2(\vec{v}) \text{ belonging to the RBM.}} \right] P_2(\vec{v})$$

$$= \left\langle \frac{\partial E}{\partial \alpha} \right\rangle_{\text{data}} - \left\langle \frac{\partial E}{\partial \alpha} \right\rangle_{\text{model}}$$

The second part is difficult as we would need to run a long MCMC, very computationally expensive! Solution (due to Hinton)

) Contrastive Divergence (C.D.)  $\rightarrow$  ignore converging the MC & take the average of  $\frac{\partial E}{\partial \alpha}$  w.r.t. configurations after  $k$  steps, usually  $k=1$ .

In practice we use stochastic gradient descent & update parameters based on mini-batches.

Let's code it up!