

DEEP — — LEARNING

UNITS ASSIGNMENT

—Done By

Jahnavi Darbhampulla
RA2011047010011

DEEP LEARNING UNIT 5 - ASSIGNMENT

D) What are Attention models in Neural Networks?

Attention Models, or attention mechanisms, are input processing techniques for neural networks that allows the network to focus on specific aspects of a complex input, one at a time until the entire dataset is categorized. The goal is to break down complicated tasks into smaller areas of attention that are processed eventually and sequentially.

How Attention Models work?

In broad strokes, attention is expressed as a function that maps a query and "set" of key value pairs to an output. One in which the query, keys, values and final output are all vectors. The output is then calculated as a weighted sum of the values, with the weight assigned to each value expressed by a compatibility function of the query with the corresponding key value.

In practice, attention allows neural networks to approximate the visual attention mechanism humans use. Like people processing an new scene, the model studies a certain point of an image with intense, "high resolution" focus, while perceiving the surrounding areas in "low resolution", then adjusts the focal point as the network begins to understand the scene.

While this is a powerful technique for improving computer vision, the most work so far with attention mechanisms has focused on Neural Machine Translation. Here, the meaning of a sentence is mapped into a fixed length vector, which then generates a translation based on that vector as a whole.

Decision Layer

Attention Model

Memory layer(s)

Context for each incoming time step

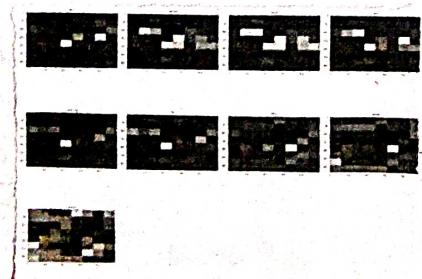
Inputs

Time

(ii) Attention over images

The Attention Mechanism is a complex ability that human beings possess. When people receive information, they can consciously ignore some of the main information while ignoring other secondary information.

In the case of Image captioning rather than compressing an entire image into a static representation, the Attention mechanism allows for salient features to dynamically come to the forefront as and when needed. This is especially important when there is a lot of clutter in the image.



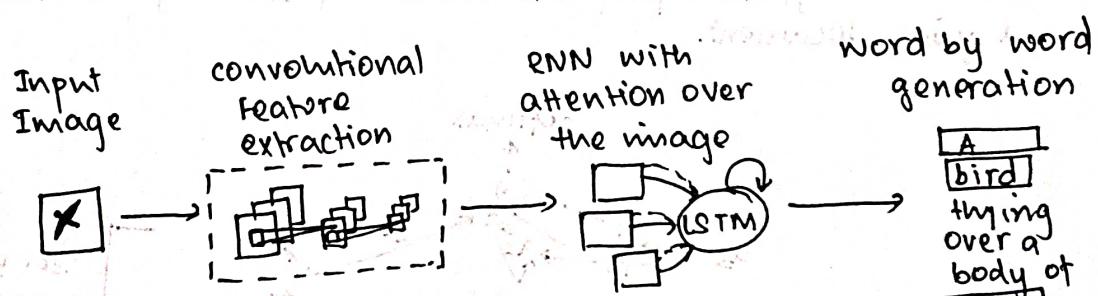
Our Aim would be to generate a caption like "two white dogs are running on the snow". To accomplish this we see how to implement a specific type of Attention mechanism called Bahdanau's Attention or local Attention.

The encoder-decoder image captioning system would encode the image, using a pre-trained convolutional Neural Network that would produce a hidden state. Then, it would decode this hidden state by using an LSTM and generate a caption.

For each sequence element, outputs from previous elements are used as inputs, in combination with new sequence data. This gives the RNN networks a sort of memory which might make captions more informative and context-aware.

But RNNs tend to be computationally expensive to train and evaluate, so in practice, memory is limited to just a few elements. Attention models can help address this problem by selecting the most relevant elements from an input image.

With an Attention Mechanism, the image is first divided into n parts, and we compute an image representation of each. When the RNN is generating a new word, the attention mechanism is focusing on the relevant part of the image, so the decoder only uses specific parts of the image.



In Bahdanau or local attention, attention is placed only on a few source positions. As Global focuses on all source side words for all target words it is computationally very expensive. To overcome this deficiency local attention chooses to focus on a small subset.

Local attention is to reduce the cost of the attention mechanism calculation. In calculation, the local attention is not to consider all the words on the source language side, but to predict the position of the source language end to be aligned at the current decoding according to a prediction function and navigate through the context window.

(ii) Hierarchical Attention

contrary to most text classification implementations, a Hierarchical Attention network (HAN) also considers the hierarchical structure of documents and includes an attention mechanism that is able to find the most important words, and sentence in a document while taking the context into consideration.

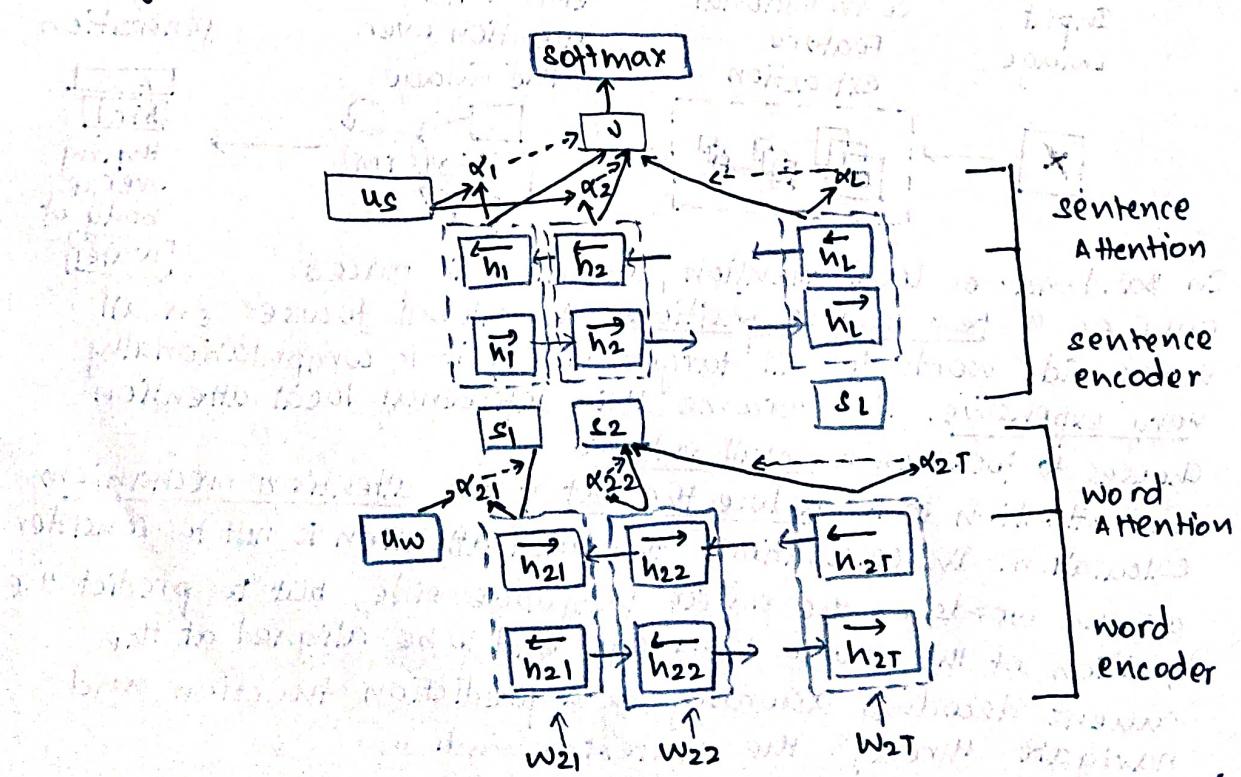
HAN finds a solution that previous works did not consider.

- Not every word in a sentence and every word in a document are equally important to understand the main message of a document.
- The changing meaning of a word depending on the context needs to be taken into consideration.

for example: The meaning of the word "pretty" can change depending on the way it is used.

"The Bouquet of flowers is pretty"] both have different meanings wrt to context.
 "The food is pretty bad"] both have different meanings wrt to context.

In this way, HAN performs better in predicting the class of a given document.



first, the network considers the hierarchical structure of documents by constructing a document representation by building representations of sentences and then aggregating those into document representation.

sentence representations are built by first encoding the word of a sentence and then applying the attention mechanism on them resulting in a sentence vector.

document representations is built in the same way, however it only receives the sentence vector of each sentence of the document as an input.

The model consists of

- the encoder which returns relevant contexts
- the attention mechanism, which computes importance weights of these contexts as one vector.

3) (i) Explain Generative Adversarial Network with a neat diagram?

Also, mention any three applications of GAN?

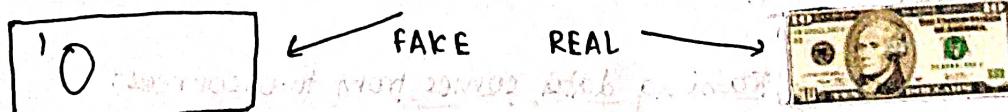
A generative adversarial Network (GAN) has 2 parts:

The generator learns to generate plausible data. The generated instances become negative training examples for the discriminator.

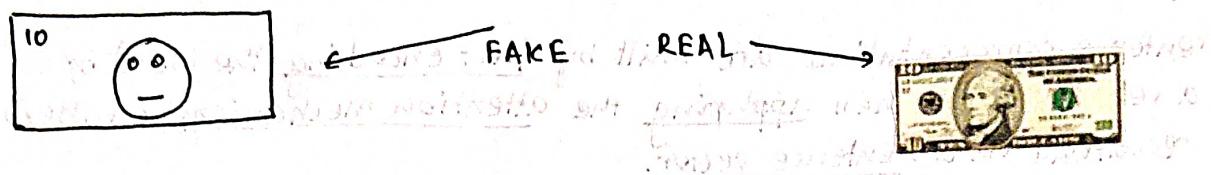
The discriminator learns to distinguish the generators' fake data from real data. The discriminator penalizes the generator for producing unplayable results.

when training begins, the generator produces obviously fake data, and the discriminator quickly learns to tell its fake.

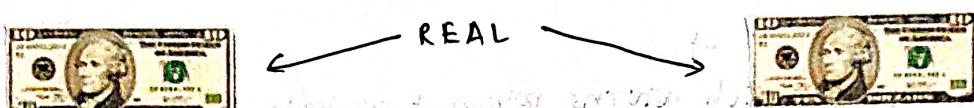
Generated Data



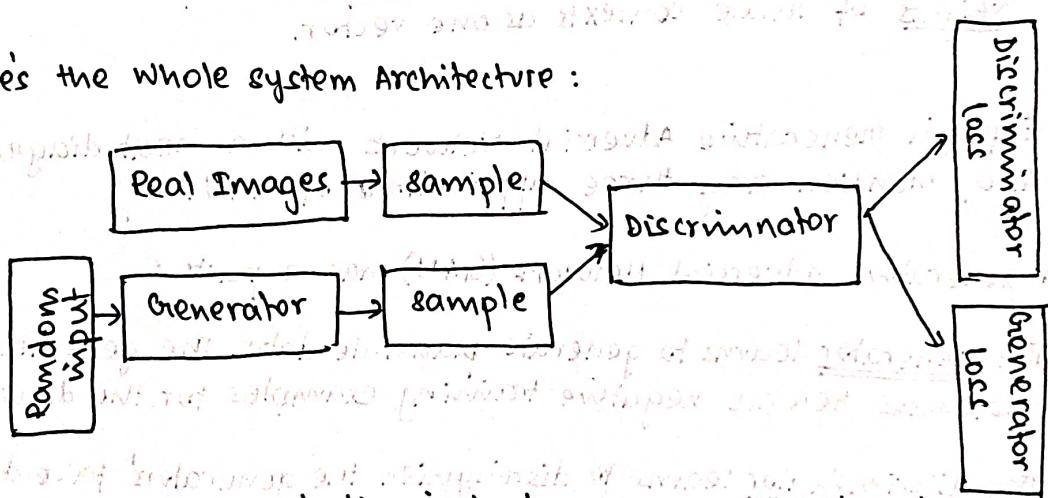
As training progresses, the generator gets closer to producing output that can fool the discriminator.



Finally if generator training goes well, the discriminator gets worse at telling the difference between real vs fake. It starts to classify fake data as real, and its accuracy decreases.



Here's the whole system architecture:



Both the generator and discriminator are neural networks. The generator output is connected directly to the discriminator input. Through backpropagation, the discriminator's classification provides a signal that the generator uses to update its weights.

Discriminator

The discriminator's training data comes from two sources:

- Real Data Instances, such as real pictures of people. The discriminator uses these instances as positive examples during training.
- Fake Data Instances, created by the generator. These are the negative examples.

generator

The generator part of GAN learns to create fake data by incorporating feedback from the discriminator. It learns to make the discriminator classify its output as real.

Generator training requires tighter integration between the generator and discriminator than discriminator training requires.

The portion of GAN that trains generator includes:

- random input
- generator network, which transforms the random input into a data instance.
- discriminator network, which classifies the generated data
- discriminator output
- generator loss, which penalizes the generator for training to fool the discriminator.

Applications of GAN

(1) Bringing Mona Lisa Back to Life

Who imagined that one day we will be able to see the expressions of the Italian noblewoman Lisa Gherardini from the famous portrait of Mona Lisa painted by the Italian artist Leonardo da Vinci. Yes, this was made possible by Deepfate.

A team of researchers at Samsung AI, Moscow used a ML system that includes few-shot adversarial learning to create expressions in the portrait. The system performs a lengthy meta-learning on a large dataset of video and after was able to frame a few one-shot learning of neural talking head models of previously unseen people.

(2) Making Photos better

With the help of GAN, a team of researchers at the top social media giant, Facebook created an approach known as (ExGANs) that can produce photo realistic personalised in-painting results that are not only perceptual but also semantically plausible. The GAN masters to fill in the missing regions of a given image while discriminator nw learns to judge the difference b/w both in-painted and real images.

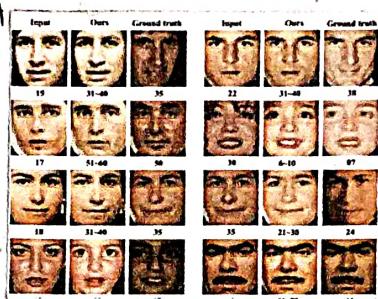
(3) Face Aging

Grigory Antipov in their 2017 titled "Face Aging with conditional Generative Adversarial Networks" use GANs to generate photographs of faces with different apparent ages from younger to older.

This is also used for deaging photographs of faces.



(2) Image Enhancement
(Making Images Better)



(1) Mona Lisa
with different
expressions

(3) Face Aging

(ii) How do you train a GAN? Explain with proper algorithm and mathematical formulations.

Algorithm

```

for number of training iterations do
    for t steps do
        • sample Minibatch of m noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ 
        • sample minibatch of m examples  $\{x^{(1)}, \dots, x^{(m)}\}$  from data generating distribution  $p_{data}(x)$ 
        • update the discriminator by ascending its stochastic gradient
    end for

```

$$\nabla_D \frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log (1 - D(G(z^{(i)})))]$$

The same procedure will be followed for training the generator by computing the loss function over a minibatch of generated images. In this step, the generator will try to fool the discriminator by generating images that look like they were sampled from the data distribution.

- sample minibatch of m noise examples/samples $\{z^1, \dots, z^m\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient.

$$\text{Update } m \text{ (gradient descent)} \\ \nabla_{\theta} \sum_{i=1}^m \log(1 - D(g(z^i)))$$

end for

The gradient based updates can use any standard gradient-descent based learning rule. We used momentum.

Algorithm Explanation

The Outer loop of the Algorithm involves iterating over steps to train the Models in the architecture. One cycle through this loop is not an epoch, it is a single update compromised of a specific batch updates to the discriminator and generator Models.

An epoch is defined as one cycle through a training dataset, where the samples in a training dataset are used to update the model weights in mini-batch.

In case of GAN, the no. of training iterations must be defined based on the size of your training dataset and batch-size.

In case of a Dataset with 100 samples, A batchsize of 10, and 500 training epochs, we would first calculate the number of batches per epoch and use this to calculate training iterations.

Next, we can use that one iteration of training results in possibly multiple updates to the discriminator and one update to generator, where number of updates to discriminator is a hyperparameter set to 1.

We therefore summarise the training algorithm with python pseudocode as follows:

```

# gan training algorithm
def train_gan(dataset, n_epochs, n_batches):
    # calculate the number of batches per epoch and global
    # batches_per_epoch = int(len(dataset) / n_batch)
    # calculate the number of training iterations
    n_steps = batches_per_epoch * n_epochs
    # gan training algorithm
    for i in range(n_steps):
        # update the discriminator model
        # ...
        # update the generator model

```

- (3) (i) What do you mean by Restricted Boltzmann Machines?
- A restricted term refers to that we are not allowed to connect the same type layer to each other. In other words, the 2 neurons of the input layer or hidden layer can't connect to each other. Although the hidden layer and visible layer can be connected to each other.
- As in this machine, there is no output layer so the question arises how we are going to identify, adjust the weights and how to measure if our prediction is accurate or not. This is where Restricted Boltzmann Machine comes to place.
- The RBM Algorithm was proposed by Geoffrey Hinton (2007), which learns probability distribution over its sample training data input. It has wide applications in different areas of supervised/unsupervised Machine learning algorithms such as feature engineering, dimensionality reduction, classification and topic modelling.

consider the movie rating recommender system, movies like Avengers, Avatar and Interstellar have strong associations with the latest fantasy and science fiction factor. Based on the user rating RBM will discover latent features that can explain the activation of movie choices.

In short, RBM describes variability among correlated variables of input dataset in terms of potentially lower numbers of unobserved variables.

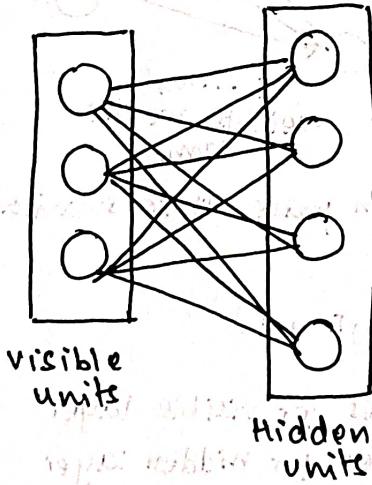
The energy function is given by,

$$E(v, h) = -a^T v - b^T h - v^T W h$$

(ii) How RBMs are used as stochastic Neural Networks?

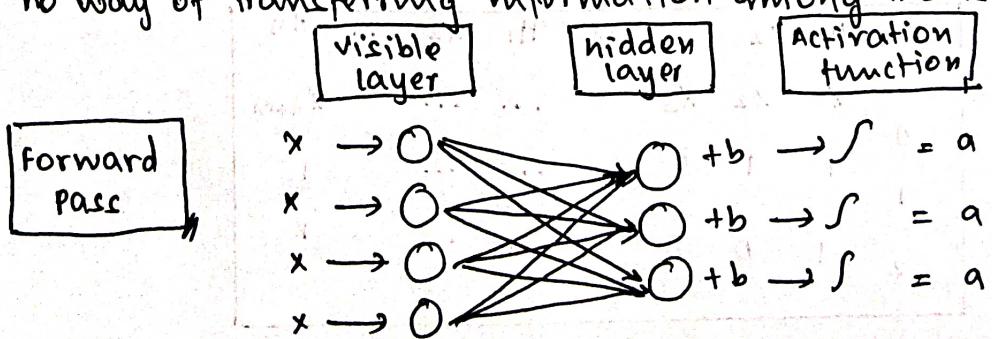
RBM is a stochastic Neural Network which means that each neuron will have some random behaviour when activated. There are 2 other layers of bias units (hidden bias and visible bias) in an RBM. This is what makes RBMs different from auto encoders.

The visible bias helps RBM to reconstruct input during backward pass.



The hidden bias RBM produce activation on forward pass and no connections between visible bias and hidden bias.

The reconstructed input is always different from the actual input as there are no connections among visible nodes and therefore no way of transferring information among themselves.



This is the first step in training an RBM with multiple inputs. The inputs are multiplied by weights and then added to bias. The result is then passed through sigmoid function.

The sigmoid function is given by

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

so the equation that we get in this step would be

$$h^{(1)} = s(v^{(0)T} w + a)$$

$h^{(1)}$ → column Matrices for hidden layer

$v^{(0)}$ → column Matrices for visible layer

These bias are new

Backward pass

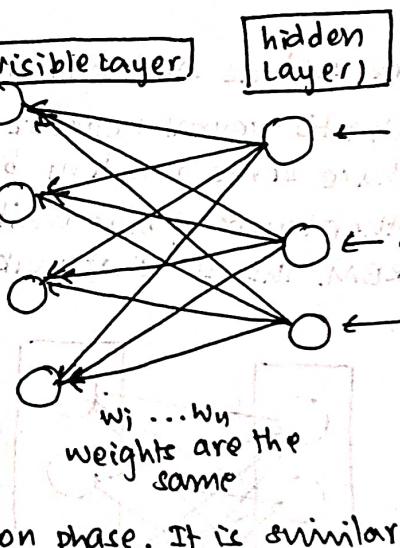
Reconstruction

reconstructions $r = b + t$

are the new output

$$r = b + t$$

and enabled GR
during MA
because no learning



Reverse phase or Reconstruction phase. It is similar to first pass but in the opposite direction.

$$v^{(1)} = s(h^{(1)} w^T + a)$$

$v^{(1)}$ → column matrices for visible layer

$h^{(1)}$ → column matrices for hidden layer

(iii) How to setup a Markov Chain Model for RBMs?

	V_1	V_2	\dots	V_m	H_1	H_2	\dots	H_n
	x_1	x_2	\dots	x_m	\dots	\dots	\dots	x_{n+m}
0	1	1	0	1
1	0	1	1	0
2	1	0	1	1

- We begin by defining our Markov chain

~~we start with a random vector $x \in \{0,1\}^{n+m}$~~

- At time step 0 we create a random vector $x \in \{0,1\}^{n+m}$.
- At time-step 1, we transition to a new value of x

⇒ How do we do this transition?

- We need to transition from a state $x = x \in \{0,1\}^{n+m}$ to $y \in \{0,1\}^{n+m}$
- sample a value $i \in \{1 \dots n+m\}$ using a distribution $q(i)$
- fix the value of all variables except x_i
- sample a new value x_i using the following conditional distribution.

• repeat the above process for many many time steps.

⇒ Where is our transition matrix? and what does it contain?

we define T such that $T_{ij} = P(x_j = y_j | x_i = x_i)$

$$P_{xy} = \begin{cases} q(i) p(y_i | x_{-i}) & \text{if } \exists i \in X \text{ so that } y_i \neq x_i \\ 0, & \text{otherwise} \end{cases}$$

where $q(i)$ is probability that x_i is the random variable whose value transitions while the value of x_{-i} remains the same.

The second term $p(y_i | x_{-i})$ essentially tells us that given the value of the remaining random variable what is the probability of x_i taking on a certain value.

The Transition Matrix T is a sparse Matrix allowing only certain transitions.

The Transition Matrix T contains the probability of transitioning from any state x to any state y .

⇒ How is it easy to create this chain? At each step we are changing only one of the n_m random variables using the following probability.

$$P(x_i = y_i | x_{-i} = x_{-i}) = \frac{p(x)}{P(x_{-i})}$$

considering the cases,

$i < m$ (we have decided to transition the value of one of the visible variables v_i to y_m)

$$\text{Then, } P(v_i = y_i | v_{-i}, h) = P(v_i = y_i | t) = \begin{cases} z & \text{if } y_i = 1 \\ 1-z & \text{if } y_i = 0. \end{cases}$$

$$\text{where } z = \sigma \left(\sum_{j=1}^m w_{ij} v_{ij} + c_i \right)$$

The above probability is easy to compute using sigmoid function.

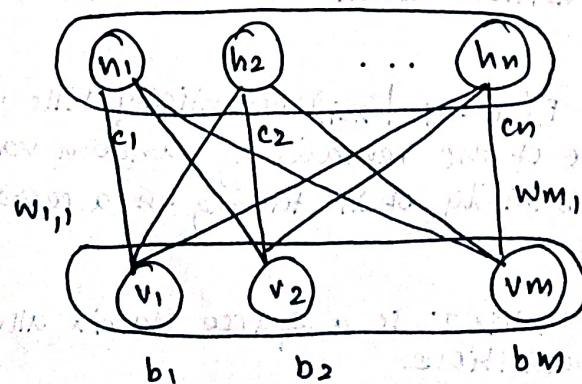
Once you compute the above probability, with probability 2 you will set the value of v_i to 1 and with probability 1-2 you will set it to 0.

so essentially at every time step you sample a v_i from uniform distribution (q_i)

And then sample a value of $v_i \in \{0, 1\}$ using the Bernoulli distribution $\text{Bernoulli}(z)$

Hence, it is easy to create this chain starting from any x_0 .

(iv) How to train RBMs using Gibbs Sampling?



$$v \in \{0, 1\}^m$$

$$E(v, h) = - \sum_i \sum_j w_{ij} v_i h_j - \sum_i b_i v_i - \sum_j c_j h_j.$$

$$\begin{aligned}\frac{\partial L(\theta)}{\partial w_{ij}} &= - \sum_H p(H|v) \frac{\partial E(v, h)}{\partial w_{ij}} + \sum_{v, h} p(v, h) \frac{\partial E(v, h)}{\partial w_{ij}} \\ &= \sum_H p(H|v) h_i v_j - \sum_{v, h} (v, h) h_i v_j \\ &= E_p(H|v) [v_i h_j] - E_{p(v, h)} [v_i h_j]\end{aligned}$$

Algorithm: RBM Training with Block Gibbs Sampling

Input: RBM ($v_1, \dots, v_m, h_1, \dots, h_n$), training batch D

Output: Learned Parameters w, b, c

init w, b, c

for all $v \in D$ do

randomly initialise $v^{(0)}$

for $t=0, \dots, k, k+1, \dots, k+r$ do

for $i=1 \dots n$ do

sample $h_i^{(t)} \sim p(h_i | v^{(t)})$

end

for $j=1 \dots m$ do

sample $v_j^{(t+1)} \sim p(v_j | h^{(t)})$

end

end

$$w \leftarrow w + n \nabla_w L(\theta) [\sigma(w v^{(t)} + c) v^{(t)} - \frac{1}{r} \sum_{t=k+1}^{k+r} (\sigma(w v^{(t)} + c) v^{(t)})^T]$$

$$b \leftarrow b + n \nabla_b L(\theta) [v^{(t)} - \frac{1}{r} \sum_{t=k+1}^{k+r} v^{(t)}]$$

$$c \leftarrow c + n \nabla_c L(\theta) [\sigma(w v^{(t)} + c) - \frac{1}{r} \sum_{t=k+1}^{k+r} \sigma(w v^{(t)} + c)]$$

end

