#### **Deep Generative Models**

#### Chapter 3: Generation by Explicit Distribution Learning

#### Ali Bereyhi

#### ali.bereyhi@utoronto.ca

#### Department of Electrical and Computer Engineering University of Toronto

#### Summer 2025

1/43

# Distribution Model: Requirements

We need to a model that computes distribution at any sample

- + Can't you just use a Softmax?!
- Well! Not if I think about the whole *x* 

  - → What if we have a continuous data space?

#### **Properties of Distributions**

A distribution model  $P_{\mathbf{w}}(x)$  should

- **1** be non-negative at any point, i.e.,  $P_{\mathbf{w}}(x) \ge 0$
- **2** add up to one, i.e.,

$$\sum_{x \in \mathbb{X}} P_{\mathbf{w}}(x) = 1 \qquad \qquad \int_{\mathbb{X}} P_{\mathbf{w}}(x) \, \mathrm{d}x = 1$$

2/43

# **Boltzmann Distribution**

Say we have a model  $f_{\mathbf{w}}\left(x
ight): \mathbb{X} \mapsto \mathbb{R}$ : we can build a distribution out of it by

**1** computing an exponential function of it: for any  $\theta$ 

(

$$\exp\left\{-\frac{f_{\mathbf{w}}\left(x\right)}{\theta}\right\}$$

**2** normalizing it to its marginalization

$$P_{\mathbf{w}}(x) = \frac{1}{\mathcal{Z}(\mathbf{w}, \theta)} \exp\left\{-\frac{f_{\mathbf{w}}(x)}{\theta}\right\}$$

where  $\mathcal{Z}\left(\mathbf{w}, \theta\right)$  is defined as

$$\mathcal{Z}\left(\mathbf{w},\theta\right) = \sum_{x} \exp\left\{-\frac{f_{\mathbf{w}}\left(x\right)}{\theta}\right\}$$

### Boltzmann's Recipe for Building Distributions

- + Can we confirm that this is a distribution?
- Sure! Just check the requirements

It is easy to see that both requirements are satisfied by this expression

1 The exponential function is always non-negative

$$P_{\mathbf{w}}(x) = \frac{1}{\mathcal{Z}(\mathbf{w}, \theta)} \exp\left\{-\frac{f_{\mathbf{w}}(x)}{\theta}\right\} > 0$$

2 The distribution adds up to one

$$\sum_{x} P_{\mathbf{w}}(x) = \sum_{x} \frac{1}{\mathcal{Z}(\mathbf{w}, \theta)} \exp\left\{-\frac{f_{\mathbf{w}}(x)}{\theta}\right\}$$
$$= \frac{1}{\mathcal{Z}(\mathbf{w}, \theta)} \underbrace{\left(\sum_{x} \exp\left\{-\frac{f_{\mathbf{w}}(x)}{\theta}\right\}\right)}_{\mathcal{Z}(\mathbf{w}, \theta)} = 1$$

4/43

#### **Building Distribution**

# **Boltzmann Distribution**

#### **Boltzmann Distribution**

The following distribution is called Boltzmann distribution at temperature heta

$$P_{\mathbf{w}}(x) = \frac{1}{\mathcal{Z}(\mathbf{w},\theta)} \exp\left\{-\frac{f_{\mathbf{w}}(x)}{\theta}\right\}$$

#### Good to Know

Boltzmann distribution is the solution to the second law of thermodynamics: it describes the distribution of micro-state at its thermal equilibrium when

- The energy of the system is described by Hamiltonian  $f_{\mathbf{w}}(x)$
- The system is isolated from the outside world 0

# Boltzmann Distribution: Components

Despite its origins in physics: for us

Boltzmann builds distribution from an arbitrary function  $f_{\mathbf{w}}$ 

We focus on the case with  $\theta = 1, 1$  and refer to

- 1 the model  $f_{\mathbf{w}}(x)$  as the energy model
- 2 the normalization term  $\mathcal{Z}(\mathbf{w})$  as the partition function

$$\mathcal{Z}\left(\mathbf{w}\right) = \sum_{x} \exp\left\{-f_{\mathbf{w}}\left(x\right)\right\}$$

#### Attention

Although the partition function does not depend on the sample point x, it still depends on the model parameter w

 $^1\mbox{We}$  don't really need to keep it general, as our model has enough parameters in  ${\bf w}$ 

Chapter 3: Explicit Methods

#### **Building Distribution**

# EBMs: Energy-based Models

Using Boltzmann distribution: we can convert any learning model  $f_{w}$  into a model for data distribution

such models are called energy-based models (EBMs)

#### **Energy-based Models**

An EBM describes a Boltzmann distribution whose energy model is given by a learnable model, e.g., a deep NN

- Is there any reason we are interested on these models?! +
- Well! They are quite generic —

# Universality of EBMs: Gaussian Example

#### Universality Argument (informal)

For almost any well-defined distribution Q, there exists a function f such that Q is expressed as a Boltzmann distribution with energy function f

**Example:** Say  $x \in \mathbb{R}$  is distributed with  $\mathcal{N}(0, 1)$ , i.e.,

$$Q\left(x\right) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2}\right\}$$

Setting the energy function to  $f(x) = x^2/2$ , it is easy to see that

$$\mathcal{Z} = \int \exp\left\{-\frac{x^2}{2}\right\} dx = \sqrt{2\pi} \rightsquigarrow P(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2}\right\}$$
$$= Q(x)$$

# Universality of EBMs: Bernoulli Example

**Example:** Say  $x \in \{0, 1\}$  is distributed with Ber(p), i.e.,

$$Q\left(x\right) = \begin{cases} p & x = 0\\ 1 - p & x = 1 \end{cases}$$

Let  $\sigma^{-1}$  denote inverse of Sigmoid function, i.e.,  $\beta = \sigma^{-1}(p)$  implies  $p = \sigma(\beta)$ . Setting energy function to  $f(x) = \sigma^{-1}(p) x$ , the partition function reads

$$\mathcal{Z} = \sum_{x} \exp\{-\sigma^{-1}(p)x\} = 1 + \exp\{-\sigma^{-1}(p)\}\$$

and thus, the Boltzmann distribution reads

$$P(0) = \frac{1}{1 + \exp\{-\sigma^{-1}(p)\}} = \sigma(\sigma^{-1}(p)) = p = Q(0)$$
  
\$\lowspace P(x) = Q(x)\$

## **Computational EBMs**

We can use Boltzmann's formula to convert

any computational model, e.g., deep NN, to a computational EBM

the computational EBMs can be roughly divided into two groups

- Boltzmann Machines (BMs)
  - → which consider a simple quadratic energy function with learnable weights
  - General BMs are the basic forms of computational EBMs
    - → Somehow like linear models for regression and classification
  - → They are inspired by the Ising or SK model in statistical physics
- Neural EBMs
  - → which consider a deep NN as the energy function
  - ↓ They are essentially the nonlinear version of BMs

10/43

#### Computational EBMs: Boltzmann Machine

#### Boltzmann Machine (BM)

Consider data  $x \in \{0, 1\}^d$ : BM models the distribution of x as a Boltzmann distribution with energy model

$$f_{\mathbf{W}}\left(x\right) = -\left\langle \mathbf{W}x;x\right\rangle$$

for a learnable  $\mathbf{W} \in \mathbb{R}^{d \times d}$  which is symmetric with diagonal zero

We can think of energy model in this case as a single linear layer



#### Few Notes on Boltzmann Machine

This form of BM is often called

#### Fully Visible BM (FVBM)

as the energy model contain only the visible data elements, i.e.,  $x_i$ 's

- + Can we have anything hidden?
- Yes! Latent representation, but we are not considering them yet

BMs can be extended to other data spaces

- Continuous data space, i.e.,  $x_i \in \mathbb{R}$ 
  - → BM describes a Gaussian distribution with learnable covariance
  - → Some people call this model the Gaussian BM
- Discrete but not binary data space, e.g.,  $x_i \in \{0, \dots, 255\}$ 
  - L→ Not too different from the basic BM

# Computational EBMs: Neural EBMs

Neural EBMs are the general form of an EBM

- Data is not restricted to have binary entries
- Energy model does not be the simple quadratic one

We can think of neural EBMs as nonlinear BMs, i.e.,



where  $f_{\mathbf{w}}\left(x\right): \mathbb{X} \mapsto \mathbb{R}$  can be a very deep network

# Sampling EBMs

Say we trained an EBM: we need to sample from it!

- + But, didn't you say sampling from a joint distribution is hard?
- Indeed! And, we can see it more clearly here!

Say we are given trained energy model vor we should compute distribution



It's exponentially hard to compute EBMs from their energy models!

Deep Generative Models

#### Training EBMs

- + Oh! What about training EBMs then? Sounds to be the same!
- Right! It's the same!

Say we are to train an EBM with energy model  $f_{\mathbf{w}}\left(x
ight)$  on dataset

$$\mathbb{D} = \left\{ x^j : j = 1, \dots, n \right\}$$

We use MLE *vert let's look at the log-likelihood* 

$$\log \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{j} \log P_{\mathbf{w}} \left( x^{j} \right) = \frac{1}{n} \sum_{j} \log \frac{\exp\left\{-f_{\mathbf{w}} \left( x^{j} \right)\right\}}{\mathcal{Z}(\mathbf{w})}$$
$$= \frac{1}{n} \sum_{j} \left[-f_{\mathbf{w}} \left( x^{j} \right) - \log \mathcal{Z}(\mathbf{w})\right]$$
$$= -\log \mathcal{Z}(\mathbf{w}) - \frac{1}{n} \sum_{j} f_{\mathbf{w}} \left( x^{j} \right)$$

#### **Training EBMs**

We can define the empirical expectation of the energy model as

$$\hat{\mathbb{E}}_{x \sim P} \left\{ f_{\mathbf{w}} \left( x \right) \right\} = \frac{1}{n} \sum_{j} f_{\mathbf{w}} \left( x^{j} \right)$$

which based on LLN goes to expectation on P as  $n \to \infty$ , since  $x^j \sim P$ 

The log-likelihood can then be written as

$$\log \mathcal{L}(\mathbf{w}) = -\left[\hat{\mathbb{E}}_{x \sim P} \left\{ f_{\mathbf{w}}(x) \right\} + \log \mathcal{Z}(\mathbf{w}) \right]$$

So, we can apply MLE training by minimizing negative log-likelihood

$$\hat{R}(\mathbf{w}) = \hat{\mathbb{E}}_{x \sim P} \{ f_{\mathbf{w}}(x) \} + \log \mathcal{Z}(\mathbf{w})$$
computed on dataset # of operation  $\alpha C^{d}$ 

# EBMs: Complexity Challenge

#### Complexity of EBM: Training and Generation

It is computationally complex to sample and train EBMs, since both

- evaluation of  $P_{\mathbf{w}}$  for sampling, and
- computation of MLE objective  $\hat{R}(\mathbf{w})$  for training

need computation of partition function which is exponentially hard

- + We saw before that sampling is expoentially hard! So, this means that we are facing two exponentially hard tasks in EBMs! Right?!
- Well! Indeed, these two hard tasks are the same thing
  - $\downarrow$  If we could sample  $P_{\mathbf{w}}$ , we could compute MLE objective as well!
- + How does it come?
- Let's take a look at training again!

# Training EBMs by Sampling Them

#### Assumption: Genie-aided Sampling

Assume that we have access to a genie which can sample  $P_{\mathbf{w}}(x)$  by only knowing its energy function  $f_{\mathbf{w}}(x)$ 

For training: we start with some  $\mathbf{w}$  and iteratively update using  $\nabla \hat{R}(\mathbf{w})$  $\downarrow$  We need  $\nabla \hat{R}(\mathbf{w})$  to compute the gradient

Let's compute the gradient

$$\nabla \hat{R}(\mathbf{w}) = \nabla \hat{\mathbb{E}}_{x \sim P} \{ f_{\mathbf{w}}(x) \} + \nabla \log \mathcal{Z}(\mathbf{w})$$
$$= \hat{\mathbb{E}}_{x \sim P} \{ \nabla f_{\mathbf{w}}(x) \} + \nabla \log \mathcal{Z}(\mathbf{w})$$
Backpropagation on energy model

So, the main challenge is to compute  $\nabla \log \mathcal{Z}(\mathbf{w})!$ 

# Training EBMs by Sampling Them

Let's focus on the challenging gradient

$$abla \log \mathcal{Z}\left(\mathbf{w}
ight) = rac{
abla \mathcal{Z}\left(\mathbf{w}
ight)}{\mathcal{Z}\left(\mathbf{w}
ight)}$$

From the definition of partition function, we have

$$\nabla \mathcal{Z} (\mathbf{w}) = \sum_{x} \nabla \exp \{-f_{\mathbf{w}} (x)\}$$
$$= -\sum_{x} \nabla f_{\mathbf{w}} (x) \exp \{-f_{\mathbf{w}} (x)\}$$
$$P_{\mathbf{w}} (x)$$

So, we could say

$$\nabla \log \mathcal{Z}(\mathbf{w}) = -\sum_{x} \nabla f_{\mathbf{w}}(x) \frac{\exp\left\{-f_{\mathbf{w}}(x)\right\}}{\mathcal{Z}(\mathbf{w})}$$

1

## Training EBMs by Sampling Them

This means that

$$\nabla \log \mathcal{Z}(\mathbf{w}) = -\sum_{x} \nabla f_{\mathbf{w}}(x) P_{\mathbf{w}}(x) = -\mathbb{E}_{x \sim P_{\mathbf{w}}} \{\nabla f_{\mathbf{w}}(x)\}\$$
$$= -\nabla \mathbb{E}_{x \sim P_{\mathbf{w}}} \{f_{\mathbf{w}}(x)\}\$$

Since the genie can give us samples from  $P_{\rm w},$  we can estimate this expectation using samples drawn from the EBM, i.e.,

$$\mathbb{E}_{x \sim P_{\mathbf{w}}} \left\{ f_{\mathbf{w}} \left( x \right) \right\} \approx \hat{\mathbb{E}}_{x \sim P_{\mathbf{w}}} \left\{ f_{\mathbf{w}} \left( x \right) \right\}$$

and hence estimate  $\nabla \log \mathcal{Z}\left(\mathbf{w}\right)$  as

$$\nabla \log \mathcal{Z}(\mathbf{w}) = -\nabla \hat{\mathbb{E}}_{x \sim P_{\mathbf{w}}} \left\{ f_{\mathbf{w}}(x) \right\}$$

20/43

# Training Metric in EBMs

So, the gradient of the MLE objective can be estimated as

$$\nabla \hat{R}(\mathbf{w}) = \nabla \hat{\mathbb{E}}_{x \sim P} \left\{ f_{\mathbf{w}}(x) \right\} - \nabla \hat{\mathbb{E}}_{x \sim P_{\mathbf{w}}} \left\{ f_{\mathbf{w}}(x) \right\}$$

#### MLE Learning for EBMs

We can perform MLE learning using the gradients of the divergence metric

$$\hat{D}(\mathbf{w}) = \underbrace{\mathbb{\hat{E}}_{x \sim P} \left\{ f_{\mathbf{w}}(x) \right\}}_{\mathsf{F}_{\mathbf{w}}} - \underbrace{\mathbb{\hat{E}}_{x \sim P_{\mathbf{w}}} \left\{ f_{\mathbf{w}}(x) \right\}}_{\mathsf{F}_{\mathbf{w}}}$$

This is an interesting expression: in MLE we try to simultaneously

- reduce the energy at data samples, i.e.,  $f_{\mathbf{w}}\left(x^{j}\right)$  for  $x^{j} \sim P$
- increase the energy at model samples, i.e.,  $f_{f w}\left( ilde{x}^{j}
  ight)$  for  $ilde{x}^{j}\sim P_{f w}$

# Sampling and Training EBMs: Summary

To train and sample with EBMs, we only need a sampling genie

Sampling Genie

It can sample  $P_{\mathbf{w}}\left(x\right)$  only knowing its energy function  $f_{\mathbf{w}}\left(x\right)$ 

Train\_EBM(D:dataset):

- 1: Initiate the energy model  $f_{\mathbf{w}}$  with some  $\mathbf{w}$
- 2: for multiple epochs do
- 3: Sample a batch of data samples  $\{x^j : j = 1, ..., n\}$  from  $\mathbb D$
- 4: Sample a batch of model samples  $\{\tilde{x}^j : j = 1, ..., n\} \leftarrow \text{SamplingGenie}(f_w)$
- 5: for j = 1, ..., n do
- 6: Compute  $\nabla f_{\mathbf{w}}(x^{j})$  and  $\nabla f_{\mathbf{w}}(\tilde{x}^{j})$  by backpropagation over energy model
- 7: end for
- 8: Update w using Opt\_avg  $\{\nabla f_{\mathbf{w}}(x^j) \nabla f_{\mathbf{w}}(\tilde{x}^j)\}$

9: end for

10: return trained energy model  $f_w$ 

# Sampling and Training EBMs: Summary

To train and sample with EBMs, we only need a sampling genie

Sampling Genie

It can sample  $P_{\mathbf{w}}\left(x
ight)$  only knowing its energy function  $f_{\mathbf{w}}\left(x
ight)$ 

Sample\_EBM():

1: Generate a sample as  $\tilde{x} \leftarrow \text{SamplingGenie}(f_w)$ 

2: return  $\tilde{x}$ 

# Key Challenge: Sampling EBMs

We need to figure out a way to sample EBMs

- + But, you said in the first section that this is exponentially hard!
- Well! We talked about direct sampling
  - └→ In direct sampling, we generate an exact sample in one shot

Since direct sampling is not tractable

we need to come up with an approximate way of sampling

A well-established way is to use Markov Chain Monte Carlo (MCMC) algorithms

### Sampling via Markov Chain Monte Carlo

Markov Chain Monte Carlo (MCMC)

In MCMC sampling, a Markov chain

$$x^{(0)} \to x^{(1)} \to \dots \to x^{(t)} \to \dots$$

is defined which starting from an easy-to-sample  $x^{(0)} \sim P^0$ , it progresses in time such that the chain converges in distribution to the target  $P_{\mathbf{w}}$ 

**Example:** Say  $x^{(0)} \sim P^0$  for an arbitrary  $P^0$ , we can define MCMC

$$x^{(t)} = \sqrt{1 - \varepsilon} x^{(t-1)} + \sqrt{\varepsilon} \mathbf{z_t}$$

for some small  $\varepsilon$  and  $z_t$  sampled i.i.d. from  $\mathcal{N}(0,1)$  at each t

 $\downarrow$  We can show that as  $t \to \infty$ , sample  $x^{(t)}$  converges to  $\tilde{x} \sim \mathcal{N}(0, 1)$ 

# MCMC: Zero and First Order Methods

MCMC methods can be roughly classified as

- Zero-order MCMC methods, which use target distribution  $P_{\mathbf{w}}(x)$  directly to build the Markov chain
  - → Simplest MCMCs in terms in terms of computation
  - $\$  Take often long time to mix, i.e., start to converge at large t
  - → Most famous over is Gibbs sampling which we learn next
- First-order methods, which use  $abla_x \log P_{\mathbf{w}}(x)$  to build Markov chain
  - ↓ More computation is needed
  - $\, \, \downarrow \, \,$  Mix faster than zero-order methods, i.e., start to converge at smaller t
  - → Most famous over is Langevin dynamics which we learn
- Second-order methods that use Hessian

  - → We do not consider them here

# MCMC<sub>0</sub>: Gibbs Sampling

Gibbs sampling is a famous zero-order MCMC algorithm



- + What is burn-in period?
- This is time needed for  $x^{(t)}$  to start converging  $P_{\mathbf{w}}$  in some sense
  - └→ The amount of time needed to forget **bad** initialization

Note that at each t we loop over dimensions  $\cdots \rightarrow d$  samplings

# MCMC<sub>0</sub>: Gibbs Sampling

#### Gibbs Sampling: Theoretical Guarantee

Under some mild conditions, sample  $x^{(T)}$  converges to sample  $\tilde{x} \sim P_{\mathbf{w}}(x)$  in distribution as  $T \to \infty$ 

- + Sounds nice! But, we still need to work with target distribution  $P_w$ !
- Yes! But we need to only work with 1D conditionals
  - └→ Similar to what we saw in AR modeling, they are much easier to sample
  - → With EBMs vor We don't need to compute partition function

#### Good to Know

Gibbs sampling is efficient for simple EBMs like basic BM

Let's consider the example of Boltzmann Machine: in this model, the data space is  $X \in \{0, 1\}^d$  and the energy model is set to

$$f_{\mathbf{W}}(x) = -\langle \mathbf{W}x; x \rangle = -\sum_{i} \sum_{\ell \neq i} W_{i,\ell} x_i x_\ell$$

where  $W_{i,\ell}$  is the entry at row i and column  $\ell$  of  $\mathbf{W}$ 

- ${\bf l}_{{\bf k}}$  W is symmetric, i.e.,  $W_{i,\ell} = W_{\ell,i}$
- $\lor$  W has zero diagonal, i.e.,  $W_{i,i} = 0$
- $\, {\scriptstyle \, {\scriptstyle \, \smile } \, } \, x_i \in \{0,1\}$  for any  $i=1,\ldots,d$

The corresponding EBM is then written as

$$P_{\mathbf{W}}(x) = \frac{\exp\left\{-f_{\mathbf{W}}(x)\right\}}{\mathcal{Z}(\mathbf{W})}$$

For Gibbs sampling, we need conditional distributions: we use Bayes rule

$$P_{\mathbf{W}}\left(\boldsymbol{x_{i}}|\boldsymbol{x_{< i}}, \boldsymbol{x_{> i}}\right) = \frac{P_{\mathbf{W}}\left(\boldsymbol{x_{< i}}, \boldsymbol{x_{i}}, \boldsymbol{x_{> i}}\right)}{P_{\mathbf{W}}\left(\boldsymbol{x_{< i}}, \boldsymbol{x_{> i}}\right)} = \frac{P_{\mathbf{W}}\left(\boldsymbol{x}\right)}{P_{\mathbf{W}}\left(\boldsymbol{x_{< i}}, \boldsymbol{x_{> i}}\right)}$$

By marginalization, we can write

$$P_{\mathbf{W}}(x_{i}) = \sum_{x_i} P_{\mathbf{W}}(x_{i})$$
$$= P_{\mathbf{W}}(x_{i}) + P_{\mathbf{W}}(x_{i})$$

So, we have

$$P_{\mathbf{W}}\left(x_{i}|x_{i}\right) = \frac{P_{\mathbf{W}}\left(x\right)}{P_{\mathbf{W}}\left(x_{i}\right) + P_{\mathbf{W}}\left(x_{i}\right)}$$

We can further simplify as

$$P_{\mathbf{W}}\left(x_{i}|x_{i}\right) = \frac{\frac{\exp\left\{-f_{\mathbf{W}}\left(x\right)\right\}}{\mathcal{Z}\left(\mathbf{W}\right)}}{\frac{\exp\left\{-f_{\mathbf{W}}\left(x_{i}\right)\right\}}{\mathcal{Z}\left(\mathbf{W}\right)} + \frac{\exp\left\{-f_{\mathbf{W}}\left(x_{i}\right)\right\}}{\mathcal{Z}\left(\mathbf{W}\right)}}$$
$$= \frac{\exp\left\{-f_{\mathbf{W}}\left(x\right)\right\}}{\exp\left\{-f_{\mathbf{W}}\left(x_{i}\right)\right\} + \exp\left\{-f_{\mathbf{W}}\left(x_{i}\right)\right\}}$$

which does not depend on the partition function!

This is indeed a simple Bernoulli distribution

$$P_{\mathbf{W}}(\mathbf{0}|x_{i}) = \frac{\exp\left\{-f_{\mathbf{W}}\left(x_{i}\right)\right\}}{\exp\left\{-f_{\mathbf{W}}\left(x_{i}\right)\right\} + \exp\left\{-f_{\mathbf{W}}\left(x_{i}\right)\right\}}$$
$$= \frac{1}{1 + \exp\left\{f_{\mathbf{W}}\left(x_{i}\right) - f_{\mathbf{W}}\left(x_{i}\right)\right\}}$$
$$= \sigma\left(f_{\mathbf{W}}\left(x_{i}\right) - f_{\mathbf{W}}\left(x_{i}\right)\right)$$
$$= \sigma\left(\Delta f_{\mathbf{W}}\left(x_{i}\right)\right)$$

So we could write

$$P_{\mathbf{W}}\left(\mathbf{x}_{i}|\mathbf{x}_{< i}, x_{> i}\right) = \begin{cases} \sigma\left(\Delta f_{\mathbf{W}}\left(\mathbf{x}_{< i}, x_{> i}\right)\right) & \mathbf{x}_{i} = 0\\ 1 - \sigma\left(\Delta f_{\mathbf{W}}\left(\mathbf{x}_{< i}, x_{> i}\right)\right) & \mathbf{x}_{i} = 1 \end{cases}$$

# Sampling BMs by Gibbs Sampling

In Boltzmann machine: the energy difference is rather easy

$$\Delta f_{\mathbf{W}}\left(x_{\langle i, x \rangle i}\right) = -\sum_{\ell \neq i} W_{i,\ell} x_{\ell} = -\langle \mathbf{w}_{i,\langle i}; x_{\langle i} \rangle - \langle \mathbf{w}_{i,\langle i}; x_{\langle i} \rangle$$

which is a linear transform on all  $x_{\ell}$  except  $x_i$ 

 $\begin{array}{l} \underline{\text{Gibbs}\_\text{Sampling}(T, \mathbf{W}):} \\ \hline 1: \ \textit{Initiate sample } x^{(0)} \\ 2: \ \textit{for } t = 1, \dots, T \ \textit{do} \\ 3: \quad \textit{for } i = 1, \dots, d \ \textit{do} \\ 4: \qquad \text{Sample } x_i^{(t)} \sim \text{Ber} \left( -\left\langle \mathbf{w}_{i, < i}; x_{< i}^{(t)} \right\rangle - \left\langle \mathbf{w}_{i, > i}; x_{> i}^{(t-1)} \right\rangle \right) \\ 5: \quad \textit{end for} \\ 6: \ \textit{end for} \\ 7: \ \textit{return } x^{(T)} \ \textit{for } T \ \textit{larger than burn-in period} \end{array}$ 

#### Training BMs by Gibbs Sampling



 $Sample_BM(T, W):$ 

- 1: Generate a sample as  $\tilde{x} \leftarrow \text{Gibbs}_\text{Sampling}(T, \mathbf{W})$
- 2: return  $\tilde{x}$

# MCMC<sub>1</sub>: Langevin Dynamics

For neural EBMs, it's better to deploy first order MCMC algorithms

zero-order MCMC algorithms often take to long to converge

A well-known first-order algorithm is Langevin dynamics

Langevin\_Sampling():

1: Initiate sample  $x^{(0)}$  and choose a converging series  $\{\epsilon_t\}$ 

2: for t = 1, ..., T do 3: Sample  $\eta_t \sim \mathcal{N}(0, 1)$ 4: Update  $x^{(t)} \leftarrow x^{(t-1)} + \frac{\epsilon_t}{2} \nabla_x \log P_{\mathbf{w}}\left(x^{(t-1)}\right) + \sqrt{\epsilon_t} \eta_t$ 5: end for

6: return  $x^{(T)}$  for T larger than burn-in period

#### Langevin Dynamics with EBMs

In EBMs, Langevin dynamics find an favorable form: we have

$$P_{\mathbf{w}}(x) = \frac{\exp\left\{-f_{\mathbf{w}}(x)\right\}}{\mathcal{Z}(\mathbf{w})}$$

Thus, we can write

$$\log P_{\mathbf{w}}(x) = -f_{\mathbf{w}}(x) - \log \mathcal{Z}(\mathbf{w})$$

Noting that  $\mathcal{Z}(\mathbf{w})$  does not depend on x, we can write

$$abla_{x}\log P_{\mathbf{w}}\left(x\right) = -
abla_{x}f_{\mathbf{w}}\left(x\right)$$

### MCMC<sub>1</sub>: Langevin Dynamics in EBMs

 $\begin{array}{l} \underline{\text{Langevin}\_\text{Sampling}():} \\ \hline 1: \ \underline{\text{Initiate sample } x^{(0)} \ \text{and choose a converging series } \{\epsilon_t\} \\ 2: \ \mathbf{for} \ t = 1, \ldots, T \ \mathbf{do} \\ 3: \quad \text{Sample } \eta_t \sim \mathcal{N}(0, 1) \\ 4: \quad \text{Update } x^{(t)} \leftarrow x^{(t-1)} - \frac{\epsilon_t}{2} \nabla_x f_{\mathbf{w}} \left(x^{(t-1)}\right) + \sqrt{\epsilon_t} \eta_t \\ 5: \ \mathbf{end for} \\ 6: \ \mathbf{return} \ x^{(T)} \ \mathbf{for} \ T \ \textit{larger than burn-in period} \end{array}$ 

This algorithm can sample from  $P_{\mathbf{w}}$  by knowing only energy function

- L→ This is indeed the sampling genie
- $\, \, {\scriptstyle \, \smile} \,$  The algorithm however still needs large T to converge

# Conservative Sampling via Langevin Dynamics

Requiring large T to converge makes it impractical for training

→ A simple remedy is to use conservative sampling for training

#### **Conservative Sampling**

In conservative sampling, we start from data sample as  $x^{(0)}$  and apply a few steps of Langevin dynamics

$$\begin{array}{l} \underline{\operatorname{CS}(k,x^{(0)}):}\\ \hline 1: \ \textit{Choose a converging} \{\epsilon_t\} & \texttt{#sample } x^{(0)} \sim P\left(x\right) \text{ is given}\\ 2: \ \textit{for } t = 1, \ldots, k \ \textit{do}\\ 3: & \ \textit{Sample } \eta_t \sim \mathcal{N}\left(0,1\right)\\ 4: & \ \textit{Update } x^{(t)} \leftarrow x^{(t-1)} - \frac{\epsilon_t}{2} \nabla_x f_{\mathbf{w}}\left(x^{(t-1)}\right) + \sqrt{\epsilon_t} \eta_t\\ 5: \ \textit{end for}\\ 6: \ \textit{return } x^{(k)} \end{array}$$

# Train EBMs by Conservative Divergence

#### With conservative sampling we can build efficient training loop for neural EBMs

 $CDTrain\_EBM(k, \mathbb{D}:dataset):$ 

1: Initiate w

2: for multiple epochs do

3: Sample a batch of data samples  $\{x^j : j = 1, ..., n\}$  from  $\mathbb{D}$ 

4: for 
$$j = 1, ..., n$$
 do

5: Sample model as 
$$\tilde{x}^j \leftarrow CS(k, x^j)$$

6: Compute conservative divergence  $ext{CD}_{k}^{j} = \nabla f_{\mathbf{w}}\left(x^{j}\right) - \nabla f_{\mathbf{w}}\left(\tilde{x}^{j}\right)$ 

7: end for

8: Update w using Opt\_avg 
$$\{CD_k^j\}$$

9: end for

10: return trained  $\mathbf{W}$ 

#### Sample EBMs after Training Conservative Divergence

#### Attention

We cannot use conservative sampling for generation, since we do not have access to data samples anymore!

Sampling\_EBM( $f_w$ ):

1: Initiate sample  $x^{(0)} \sim \mathcal{N}(0, 1)$  and choose converging  $\{\epsilon_t\}$ 

2: for  $t = 1, \ldots, T$  do

3: Sample 
$$\eta_t \sim \mathcal{N}(0,1)$$

4: Update 
$$x^{(t)} \leftarrow x^{(t-1)} - \frac{\epsilon_t}{2} \nabla_x f_{\mathbf{w}} \left( x^{(t-1)} \right) + \sqrt{\epsilon_t} \eta_t$$

5: end for 6: return  $x^{(T)}$  for T larger than burn-in period

#### Wrap Up

In summary, we could say

- EBMs convert any model to a distribution using Boltzmann formula
- They are not easy to sample *vorthis* impacts both sampling and training
  - → We cannot sample them directly, as it is exponentially complex
  - - └→ These approaches are still slow, as they need time to converge

#### Good to Know

EBMs can also be extended to include latent variables as well

- We get to know latent space in the next section
- Since EBMs are not very practical, we leave latent-space EBMs
  - └ If interested, take a look at restricted Boltzmann machine

# Grandfather of Diffusion Models

- + If not practical, why did we spend time learning them?!
- What we learned here serves us well in understanding diffusion models

Indeed, EBMs are the grandfathers of diffusion models: Langevin dynamics was inspired by Brownian motion which is

fundamental model for describing the diffusion process

Looking at the process of sampling with Langevin dynamics, Hyvärinen came up with an interesting conclusion in 2005

We would be fine if we try to learn so-called score function  $\nabla_x \log P(x)$ 

# Score Matching: A Short Overview

Hyvärinen observed that in Langevin dynamics, we update

$$x^{(t)} \leftarrow x^{(t-1)} + \frac{\epsilon_t}{2} \nabla_x \log P_{\mathbf{w}} \left( x^{(t-1)} \right) + \sqrt{\epsilon_t} \eta_t$$

which means that we mainly make use of the score function

$$s_{\mathbf{w}}\left(x\right) = \nabla_x \log P_{\mathbf{w}}\left(x\right)$$

Hyvärinen showed that it is feasible to match this learnable function to the true score function, i.e., to find a way to feasibly minimize

$$\mathcal{J}(\mathbf{w}) = \mathbb{E}_{x \sim P} \left\{ |s_{\mathbf{w}}(x) - s(x)|^2 \right\}$$

This is called score-matching which is efficient alternative way to train EBMs

→ We learn them in Chapter 5, as diffusion models also rely on learning score!