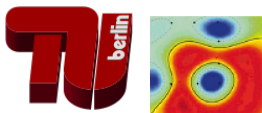
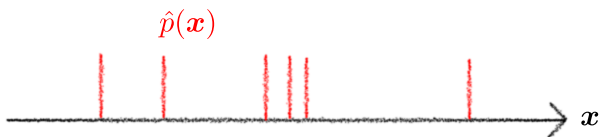


Lecture 6: Restricted Boltzmann Machines

Machine Learning 2



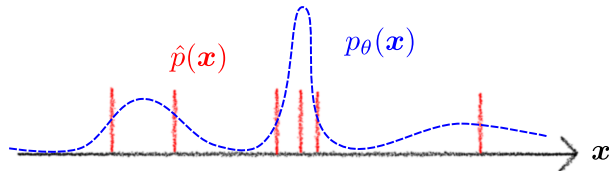
Density Estimation



Goal:

- Infer the probability distribution $p(x)$ that has been generating these few data points.

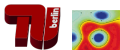
Density Estimation



- ▶ Try to learn from $\hat{p}(x)$ a distribution $p_{\theta}(x)$ that is as close as possible to the true distribution $p(x)$.
- ▶ Divergence between distribution can be measured using e.g. the KL divergence:

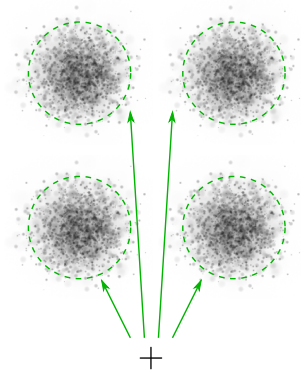
$$\text{KL}(\hat{p}||p_{\theta}) = \int_{\mathbf{x}} \hat{p}(\mathbf{x}) [\log \hat{p}(\mathbf{x}) - \log p_{\theta}(\mathbf{x})]$$

- ▶ Other divergence measures exist (e.g. α -divergence, Wasserstein distance, etc.).



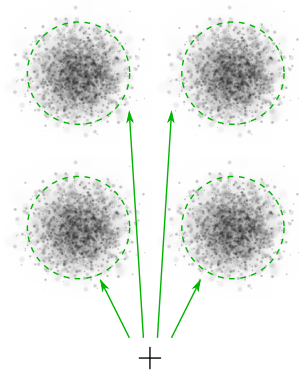
Various Approaches to Density Modeling

1. Local Mixture

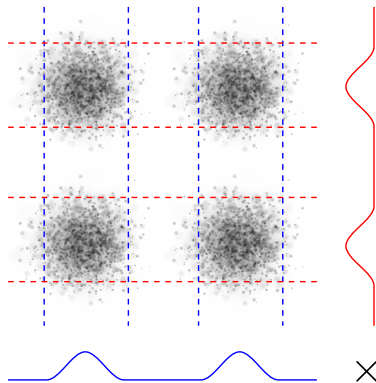


Various Approaches to Density Modeling

1. Local Mixture



2. Product of Experts



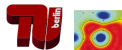
Various Approaches to Density Modeling

Gaussian Mixture Model (*studied in ML1*)

$$p_{\theta}(\mathbf{x}) = \sum_{j=1}^N \alpha_j \exp \left(\underbrace{-(\mathbf{x} - \boldsymbol{\mu}_j)^{\top} \boldsymbol{\Sigma}_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j)}_{\text{local containment}} \right)$$

Restricted Boltzmann Machine (*today's lecture*)

$$p_{\theta}(\mathbf{x}) = \frac{1}{Z} \exp(\mathbf{a}^{\top} \mathbf{x}) \prod_{j=1}^N \left(1 + \exp \left(\underbrace{\mathbf{w}_j^{\top} \mathbf{x} + b_j}_{\text{global projection}} \right) \right)$$



Restricted Boltzmann Machines

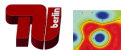
A Restricted Boltzmann machine (RBM) is a *parameterized* probability distribution over $\mathbf{x} \in \{0, 1\}^d$ defined as:

$$p_{\theta}(\mathbf{x}) = \sum_{\mathbf{h} \in \{0, 1\}^N} p_{\theta}(\mathbf{x}, \mathbf{h})$$

where

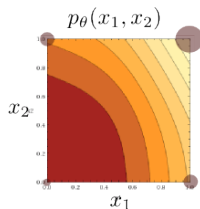
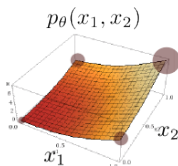
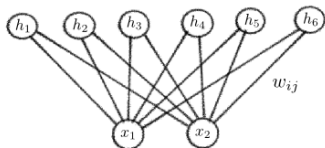
$$\begin{aligned} \blacktriangleright \quad p_{\theta}(\mathbf{x}, \mathbf{h}) &= \frac{1}{Z} \exp \left(\sum_{i,j} x_i w_{ij} h_j + \sum_i a_i x_i + \sum_j b_j h_j \right) \\ &= \frac{1}{Z} \exp \left(\underbrace{\mathbf{x}^{\top} W \mathbf{h} + \mathbf{a}^{\top} \mathbf{x} + \mathbf{b}^{\top} \mathbf{h}}_{-E_{\theta}(\mathbf{x}, \mathbf{h})} \right), \end{aligned}$$

- ▶ Z is the *partition function* that normalizes the probability distribution
- ▶ $\theta = \{W, \mathbf{a}, \mathbf{b}\}$ are the parameters of the model.



Restricted Boltzmann Machines

Example: RBM with $d = 2$ visible units and $N = 6$ hidden units.



- ▶ The RBM assigns probability mass on vertices $\{0, 1\}^2$ of the input hypercube.
- ▶ The value of the probability function at other locations is irrelevant.

The Feature View of an RBM

The probability distribution

$$p_{\theta}(\mathbf{x}) = \sum_{\mathbf{h} \in \{0,1\}^N} p_{\theta}(\mathbf{x}, \mathbf{h})$$

with

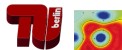
$$p_{\theta}(\mathbf{x}, \mathbf{h}) = \frac{1}{Z} \exp \left(\sum_{i,j} x_i w_{ij} h_j + \sum_i a_i x_i + \sum_j b_j h_j \right)$$

can be rewritten as:

$$p_{\theta}(\mathbf{x}) = \frac{1}{Z} \exp \left(\underbrace{\mathbf{a}^{\top} \mathbf{x} + \sum_{j=1}^N \text{softplus}(W_j^{\top} \mathbf{x} + b_j)}_{-F_{\theta}(\mathbf{x})} \right)$$

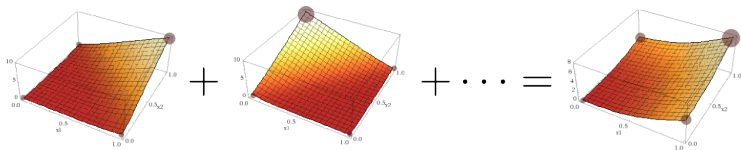
where $\text{softplus}(t) = \log(1 + e^t)$.

[Derivations on the blackboard]



RBM's can Model any Binary Distributions

Any *binary* (log-)probability distribution $\log p(x)$ over $\{0, 1\}^d$ can be written as a superposition of many softplus functions:



Remark: Universality does not hold for the *continuous* distributions in $[0, 1]^d$. Reason: a sum of convex functions can only produce a convex function. Any non-convex log-probability function cannot be represented.

Learning RBMs

- ▶ **Goal:** Learn a set of RBM parameters $\theta = \{W, \mathbf{a}, \mathbf{b}\}$ from an ensemble of data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$.
- ▶ **Optimization problem:**

$$\min_{\theta} D_{\text{KL}}(p \parallel p_{\theta})$$

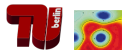
where

$$p(\mathbf{x}) = \frac{1}{n} \sum_{k=1}^n \delta(\mathbf{x} - \mathbf{x}_k)$$

and

$$p_{\theta}(\mathbf{x}) = \frac{1}{Z} \exp(-F_{\theta}(\mathbf{x}))$$

- ▶ **Idea:** Minimize the objective using gradient descent.



Learning RBMs

- ▶ **Result 1:** The gradient of the KL divergence is given by:

$$\nabla_{\theta} D_{\text{KL}}(p \parallel p_{\theta}) = \langle \nabla_{\theta} F_{\theta}(\mathbf{x}) \rangle_p - \langle \nabla_{\theta} F_{\theta}(\mathbf{x}) \rangle_{p_{\theta}}$$

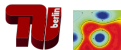
and depends on the RBM free energy $F_{\theta}(\mathbf{x})$.

[Derivations on the blackboard]

- ▶ **Result 2:** The gradient of the free energy is given by

$$\nabla_{w_{ij}} F_{\theta}(\mathbf{x}) = \underbrace{\text{sigm}(W_j^{\top} \mathbf{x} + b_j)}_{p(h_j=1|\mathbf{x})} \cdot x_i$$

[Derivations on the blackboard]



Learning RBMs

- ▶ **Problem:** Evaluating the expectation $\langle \cdot \rangle_{p_\theta}$ in the gradient requires to evaluate a function for an exponentially large number of states (2^d).
- ▶ **Idea:** Replace the expectation operator by an empirical estimate, and sample only n realizations from the $p_\theta(x)$.
- ▶ **Remark:** Ideally, the sample must be unbiased so that for n large, its statistics converge to those of the original probability distribution p_θ .
- ▶ **Question:** How do we sample from p_θ ?

Sampling in an RBM

► **Observation:**

$$p(\mathbf{h}|\mathbf{x}) = \prod_j \underbrace{p(h_j|\mathbf{x})}_{\text{Ber}(\text{sigm}(W_j^\top \mathbf{x}))} \quad \text{and} \quad p(\mathbf{x}|\mathbf{h}) = \prod_i \underbrace{p(x_i|\mathbf{h})}_{\text{Ber}(\text{sigm}(W_i^\top \mathbf{h}))}$$

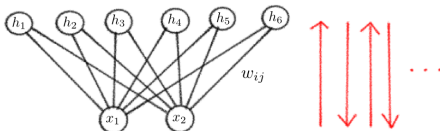
[Derivations on the blackboard]

- **Consequence:** Knowing \mathbf{x} , we can sample $\{h_j\}$'s all at once and knowing \mathbf{h} , we can sample $\{x_i\}$'s all at once.

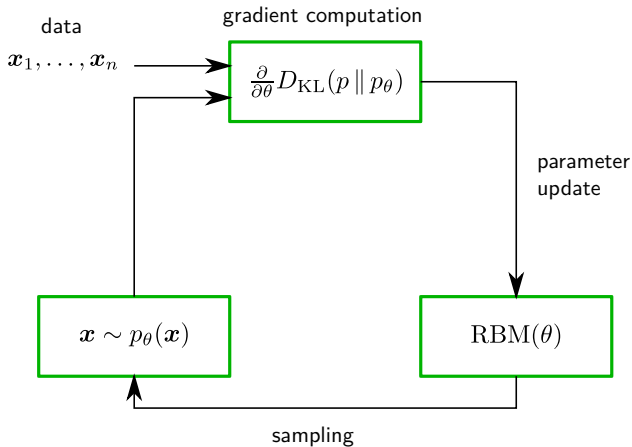
$$\mathbf{h} \leftarrow \text{Ber}(\text{sigm}(W^\top \mathbf{x} + \mathbf{b}))$$

$$\mathbf{x} \leftarrow \text{Ber}(\text{sigm}(W \mathbf{h} + \mathbf{a}))$$

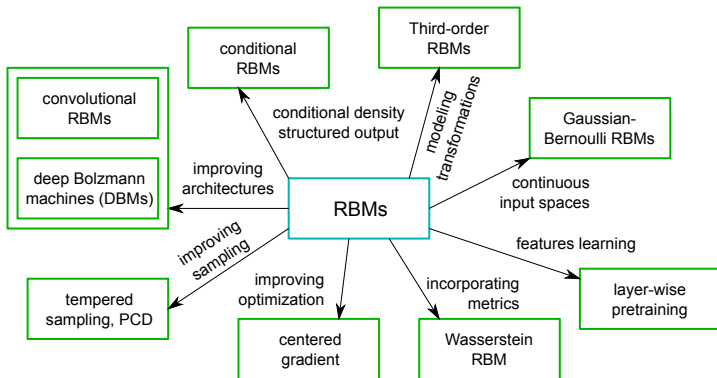
where $\text{Ber}()$ and $\text{sigm}()$ apply element-wise.



Wrapup



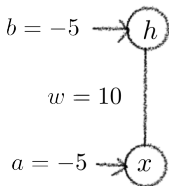
The world of Boltzmann Machines



Improving Sampling

- ▶ **Question:** How long do we need to run the Gibbs sampling procedure for it to produce an unbiased sample of the data distribution?
- ▶ **Example: Two-Unit RBM**

$$p(x) = \frac{1}{Z} \exp(xwh + bh + ax)$$



| | $h = 0$ | $h = 1$ |
|---------|---------|---------|
| $x = 0$ | 0.497 | 0.003 |
| $x = 1$ | 0.003 | 0.497 |

Will move from one mode to the other with probability 0.006 (on average every 167 iterations).

Improving Sampling

Strategies to overcome the sampling problem:

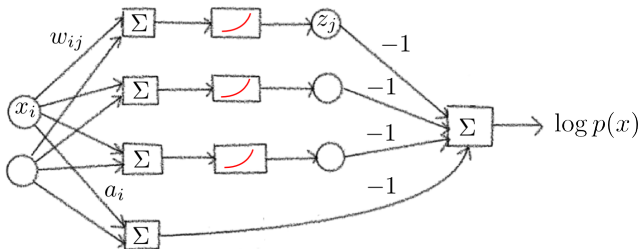
- ▶ **Annealing:** Consider a simpler model θ_{easy} and run the alternate Gibbs sampler while gradually restoring θ to its original value.
- ▶ **Improved architectures:** Design Boltzmann machine architectures in which sampling is faster.
- ▶ **Approximate sampling:**
 - ▶ *Contrastive divergence:* Sample by running a predefined number k of Gibbs updates starting from an existing data point (CD- k).
 - ▶ *Sampling by learning:* Use the impulse of learning to force the Gibbs sampler to escape from its current mode (persistent CD).

Initializing Neural Networks with RBMs

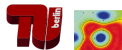
Reminder: The free energy of an RBM is given by:

$$F(\mathbf{x}) = -\mathbf{a}^\top \mathbf{x} - \sum_{j=1}^N \underbrace{\text{softplus}(W_j^\top \mathbf{x} + b_j)}_{z_j}$$

Neural Network Equivalent:



Idea: Reuse the first layer of the network to train on the discriminative problem.

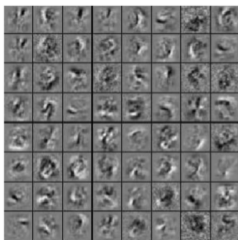


Initializing Neural Networks with RBMs

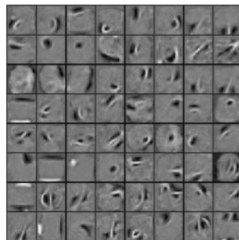
Example: MNIST



First layer weights
without pretraining



First layer weights
with pretraining



Source: Erhan, JMLR, 2010