

Boltzmann Machine and Restricted Boltzmann Machine

- Boltzmann Machine general settings : Binary variables formed undirected graph . with Energy-Based Distribution, we got hidden / latent variables and visible/obliged variables for different usage.

- Restricted Boltzmann Machine

- Restricted: No inter-connection between inter-layers. which means that the units within one layer (hidden or visible) are conditionally independent.

- Formal Representation:

$$p(v, h) = \frac{1}{Z} \exp(-\mathcal{E}(v, h)), \quad \mathcal{E}(v, h) = -b^T v - c^T h - v^T w h$$

- Conditional Distributions:

$$p(h|v) = P(v, h) / p(v) = \frac{1}{p(v)} \cdot \frac{1}{Z} \exp \left\{ b^T v + c^T h + v^T w h \right\} = \frac{1}{Z} \exp(c^T h + v^T w h)$$

For binary cases:

$$p(h_j=1|v) = \delta(c_j + v^T w_{:,j})$$

- Training RBM

we can learn explicit form of $p(h|v)$, and if we want to know $p(v)$, Sampling techniques can be used .

Deep Belief Networks

- General Structure

Top layer (Deepest Hidden Layers) is an RBM, and others (from hidden to visible) are directed.

- Generating Samples

First we sample first top layers (RBM) and then sampling data using ancestral Sampling.

- Training phase: First training an RBM to capture distribution of visible variables.

then adding hidden layers by training more RBMs for capture top hidden layers' distribution. Note that for training, we don't specify top layers must follow in a directed graph style.

- We always initialize a MLP using weights of a trained deep belief network

Deep Boltzmann Machines

- Stacking many RBMs and for visible units, the value could be real or binary, for hidden layers, the value is binary. RBM means no connection within layer. This structure can be interpreted as two separate layers with conditionally independent. (odd number and even number independent)

- Interesting Properties

No intralayer interaction makes it possible to use fix-point equations to maximize variational lower bound and find the optimal mean field expectations.

- DBM mean field inference

Because here $P(h_1, h_2 | v) \neq P(h_1 | v) P(h_2 | v)$, we need to use variational inference method, to define a Q to approximate $P(h | v)$

As other algorithms of learning manifold of training data, we can use mean-field variational inference to approximate $P(h | v)$ and consequently learning parameters W , to build up the whole DBM.

- Learning parameters of DBM

We may confront partition function when we doing parameter updating, thus training a DBM involves variational inference of $\bar{P}(h | v)$ and sampling methods.

- Layer-wise Pretraining

Layer-wise pretraining reveals most power of DBM, thus layer-wise pretraining acts as a greedy method, starting from two-layer RBM and greedy training by

adding layers.

- Jointly Training Deep Boltzmann Machines
 - Centered deep Boltzmann Machine: Reparametrize the model by $\tilde{x} = x - \mu$, creating a well conditioned Hessian.
 - MP-DBM: Change the target to obtain a true value of remaining units. adding this mechanism is equivalent to training end-to-end in deep learning context, but gradient direction and value is not a probabilistic explanation.
 - Dropout used in MP-DBM: Deletion and latent.

Boltzmann Machines for real data

Gray scale image can be used for modeling scalar of pixel values.

- Gaussian-Bernoulli RBMs.

Most intuitively, $P(v|h)$ in PBM, we define v as a continuous variable, then hidden units related to visible ones form a Gaussian Distribution together.

Corresponding Energy function: We eliminate cross-term.

- Undirected Model of Conditional Covariance.

When we considering Gaussian PBM for real values, we often care more about mean than variance, which means that we may lose information interacted with pixels (Natural images thus not OK for this)

- Mean and Covariance PBM

Part of the model for mean and part of the model for covariance. unlike VAE, here we use energy function and explicit calculating inference.

- Mean-product of Student t-distributions

- Spike and Slab Restricted Boltzmann Machines

mean is given by $(\underline{h} | \underline{OS}) W^T$ ^{binary} \rightarrow continuous, thus they together

act as gating of different densities to decide mean and co-variance

Convolutional Boltzmann Machines

Solving image's input by using convolutions; main difficulties here is not matrix multiplication, but how we deal with pooling.

1. Max-pooling: How we calculate pooling? we need compute $l \times k$ values.

2 Probabilistic max pooling: Indicating which state is most possible by a R^{ht} vector but fixed size, we cannot adjust pooling dynamically.

3. Padding: May cause fail on hidden units.

Boltzmann Machines for Structured or Sequential Outputs

RNN-RBM: Using an RNN for outputting weights by layers in RBM. and training means back-propagation through time.

Other Boltzmann Machines

Combining with some ideas of supervised learning.

Back-propagation through Random Operations

Idea: For generating samples from model, we need to add random factor on input

Reparametric trick: $y \sim N(\mu, \delta^2)$, $y = \mu + \delta z$, $z \sim N(0, 1)$

The z is independent from model, we can train μ and δ ,
using back propagation.

- Back-propagation through Discrete Stochastic Operations.

If the input is discrete, we can use $\nabla_y J(y)$ to estimate loss of different y . Monte Carlo Methods)

RBMFORCZ Algorithm:

Estimating $\nabla_w \left[\frac{\partial J(y)}{\partial w} \right]$ using monte carlo methods, also we can add $b^T w$ as bias in case unstable updating gradients.

Directed Generative Nets

- Sigmoid Belief Networks

A vector of binary state s , with each element of the state influenced by its ancestors:

$$P(s_i) = \delta(\sum_j w_{j,i} s_j + b_i)$$

A sigmoid belief network have universal estimation property.

- Differentiable Generator Networks

- Generator: A neural network which can map a random variable into a generative sample, $P(X|z)$ can represent this operation.

- Two ways of generating samples (by modeling different kind of distribution)
 - Modeling marginal distribution $P(x)$, describing all the density pattern of training data.
 - Modeling conditional distribution $P(x|z)$, where we control generative samples by modifying latent variable z .
- For generative modeling's representative ability, it can reach the ability of supervised learning as long as we choose appropriate z and $P(X|z)$.

- Variational Auto Encoder

- Compared with normal auto encoder:

Auto encoder: Capture z , which indicates manifold of x , $P(z|x)$.

VAE: Capture $q(z|x)$, variational inference of z , have generation ability.

- Compared with Generative Adversarial Network

GAN: Generative model combining with a discriminator

VAE: Generative model combining with a variational inference network

- Analysis of generating blurry images: Maximum likelihood and MLE for maximizing the ELBO.

- Generative Adversarial Networks

- A possible problem of training GAN: Generative model and discriminative model may dominate the loss in turn.
- In real usage, we sometimes change target to maximize the classification loss.
- GAN can capture features/patterns which do not exist in training data.

- Generative Moment Matching Networks

Assume that training data and samples generated by model are as similar as possible, the metric of similarity is moment: $\mathbb{E}_x \prod_i x_i^{n_i}$

- Convolutional Generative Networks

Transpose of convolutional operation $\hat{x} = w^T w \cdot x$, involves design a up-pooling method to ensure information flow back.

- Auto-Regressive Networks

Constraining dependencies - Auto-Regressive Networks: All the variables in sigmoid network are visible, we inference $P(X_i | X_{\bar{i}})$ by designing structure and calculating probability.

- Linear Auto-Regressive Networks: only stands for $P(X_i | X_j) \forall j \neq i$.

using - Neural Auto-Regressive Networks

graphical model - Can deal with discrete values by setting units corresponding to different parameter dimensions of data.

sharing - Using hidden units for gathering information groups of x .

- Neural Auto-Regressive Networks Estimator: Sharing weights from W_{ij} .

Drawing Samples from Autoencoders

- Markov chain Associated with any denoising autoencoder.

$x \rightarrow C(\tilde{x}|x) \rightarrow h = f(\tilde{x}) \rightarrow w = g(h)$ Construct a markov chain iteratively sampling encoding and decoding.
 $x = p(x|w) \leftarrow$

- Clamping and Conditional Sampling: $P(X_t | X_0)$

- Well-Back Training Procedures

Generative Stochastic Network

$P(x^{(k)} | h^{(k)}) \rightarrow P(h^{(k)} | x^{(k-1)}, h^{(k-1)})$, in one step of Markov Chain, merging $h^{(k-1)}$ and $x^{(k-1)}$ together to get $h^{(k)}$.

- Discriminative G.S.N: $P(y|x)$ for supervised learning

Other Sampling methods: Recovering diffusion process of a structured probabilistic model.

Evaluating Generative models and conclusion

- Human cannot find some potential mistakes of generative model, though data seems to be reasonable.

- we model $P(h|x)$, indicating the mani-fold of data the motif of generative modeling