

Learned Approximate Inference: Wake-Sleep Algorithm

Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics in Approximate Inference

- Task of Inference
- Intractability in Inference
 1. Inference as Optimization
 2. Expectation Maximization
 3. MAP Inference and Sparse Coding
 4. Variational Inference and Learning
 5. Learned Approximate Inference

Overview

- Inference can be thought of as optimization that increases the value of a function \mathcal{L}
- Explicitly performing optimization is expensive
- This expense can be avoided by learning to perform approximate inference
- Optimization is thought of as a function f that maps an input \mathbf{v} to an approximate distribution $\mathbf{q}^* = \arg \max_{\mathbf{q}} \mathcal{L}(\mathbf{v}, \mathbf{q})$ Recall $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, \mathbf{q})$ is lower bound on $p(\mathbf{v} ; \boldsymbol{\theta})$
- We can approximate it with a neural network that implements a neural network that is an approximation $\hat{f}(\mathbf{v}; \boldsymbol{\theta})$

Wake-Sleep Algorithm

- Difficulty of training model to infer h from v :
 - No supervised training set to train the model
 - Given a v we do not know the appropriate h
 - Mapping from v to h depends on model family, and evolves through learning process as θ changes
- Wake-sleep resolves problem by drawing samples of both h and v from model distribution
 - E.g., in directed model perform ancestral sampling beginning at h and ending at v
- Inference network can be trained to perform reverse mapping:
 - Predicting which h caused the present v

Drawback of Wake-Sleep

- We will only be able to train the inference network on values of v that have a high probability under the model
- Early in training the model does not resemble the distribution
- So the inference network does not have the opportunity to learn from samples that resemble the data

Biological Dreaming Explanations

- Two explanations of dream-sleep in human beings and animals:
 1. *Partition Function*: Dreams provide negative phase samples that Monte Carlo training algorithms use to approximate $-\nabla_{\theta} \log Z(\theta)$
 - In gradient ascent, we maximize $\nabla_{\theta} \log p(\mathbf{x}; \theta) = \nabla_{\theta} \log \tilde{p}(\mathbf{x}; \theta) - \nabla_{\theta} \log Z(\theta)$
 - Involves maximizing $\log \tilde{p}$ where data occurs and minimizing $\log Z$ where model samples occur
 2. *Wake-Sleep*: Dreams provide samples from $p(\mathbf{h}, \mathbf{v})$
 - used to train inference network to predict \mathbf{h} given \mathbf{v}
- Second explanation is more satisfying than the partition function explanation

Disadvantage of first explanation

- Monte Carlo algorithms do not perform well if they are run using only the positive phase of the gradient for several steps and then with only the negative phase of the gradient for several steps
- In contrast, human beings are awake for several consecutive hours and sleep for several consecutive hours
 - This schedule does not support Monte Carlo training of an undirected model

Learning algorithms based on maximizing \mathcal{L}

- Can be run for prolonged periods of improving q and prolonged periods of improving θ
- If the role of biological dreaming is to train networks for predicting q , then this explains how animals are able to remain awake for several hours
 - Longer they are awake, greater the gap between \mathcal{L} and $\log p(\mathbf{v})$, but \mathcal{L} will remain a lower bound
- And to remain asleep for several hours
 - Generative model itself is not modified during sleep ⁸