

Structured Probabilistic Modeling for Deep Learning

Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics in Structured PGMs for Deep Learning

0. Overview

1. Challenge of Unstructured Modeling

2. Using graphs to describe model structure

3. Sampling from graphical models

4. Advantages of structured modeling

5. Learning about Dependencies

6. Inference and Approximate Inference

7. The deep learning approach to structured probabilistic models

1. Ex: The Restricted Boltzmann machine

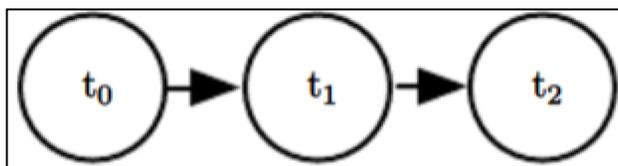
Modeling formalisms in Deep Learning

- Deep learning researchers draw on many modeling formalisms to guide their design efforts and describe their algorithms
- One of these formalisms is the idea of *Structured probabilistic models* (also known as probabilistic graphical models, *PGMs*)
- They are a key ingredient of the most important research topics in deep learning

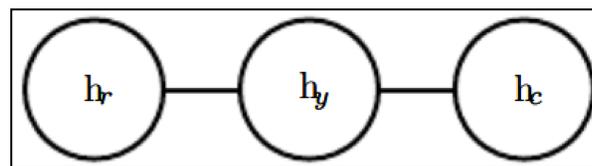
Structured Probabilistic models

- They are a way of describing probability distributions using a graph to describe which variables interact with each other directly
- Graph is used in the sense of graph theory
 - Vertices connected to one another by edges

Directed PGM



Undirected PGM



- Because structure is described by a graph they are called *graphical models* or *PGMs*

What do PGMs offer to Deep Learning?

- Graphical models community has developed:
 - Models
 - Training algorithms
 - Inference algorithms
- Some are very useful for deep learning
 - We will highlight which of these models have proven most useful

Treatment of PGMs in Deep Learning

- Deep learning differs from general PGMs in:
 1. Model structures
 2. Learning algorithms
 3. Inference procedures
- We explain differences and explain reasons

Plan of PGM Discussion

1. Challenges of large scale probabilistic models
2. Graphs to describe probability distributions
 - It has its own complications
 - E.g., which variables need to interact directly?
 - Discuss two approaches to overcome difficulty
 - By learning about dependencies
3. Unique emphasis in deep learning approaches

Challenge of Unstructured Modeling

- Goal of deep learning is to scale ML to the kind of challenges needed to solve AI
 - This means being able to understand high-dimensional data with rich structure
 - Understand natural images
 - A natural image is one captured by a camera in an ordinary environment, as opposed to a synthetically rendered image, a screenshot of a web page
 - Audio waveforms representing speech
 - Documents with multiple words, punctuation characters
- The classification task of machine learning is a limited goal

Classification is a limited goal

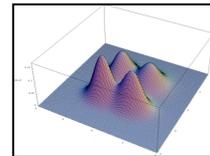
- Classification task
 - Input from a distribution is summarized with a label
 - Object in photo
 - Word in audio
 - Topic of document
 - Discards most of the input
 - Ignore background when recognizing object
 - Produces a single output
 - Or a distribution over values of that single output
- Probabilistic models can do many other tasks
 - Often more expensive than classification

Graphical models for other tasks

- Tasks
 - Produce many output values
 - Understand entire input with no ignoring of sections

- Some tasks are:

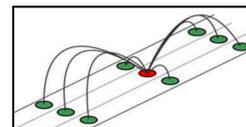
1. Density estimation



2. Denoising



3. Missing value imputation

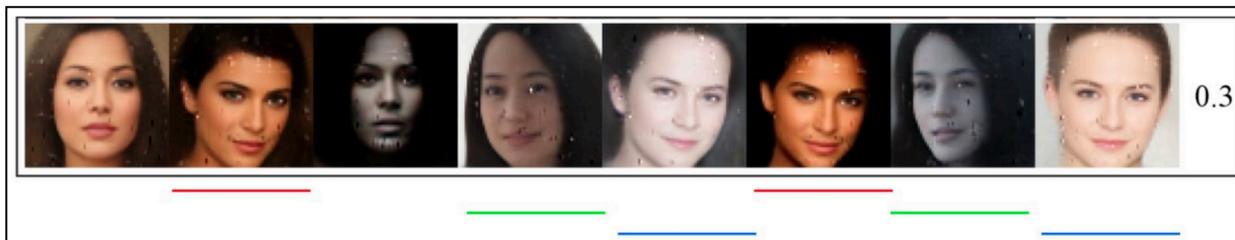
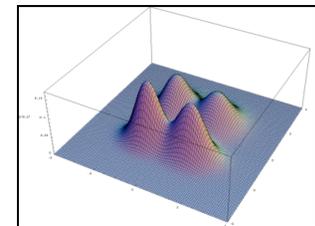


4. Sampling



1. Density Estimation Task

- Given \mathbf{x} , return $p(\mathbf{x})$ under generating distribution
 - Although one output, need entire input space
 - If element of \mathbf{x} is unusual, assign low probability
- Modes of the distribution must be preserved
 - GANs have a famous problem with this
 - Called the “Mode Collapse” problem

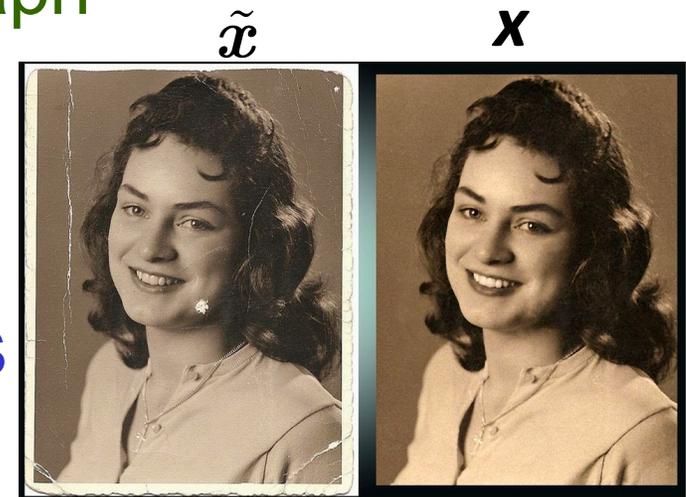


Partial collapse happens often.

Images with same underlined color look similar and the mode starts collapsing.

2. The Denoising Task

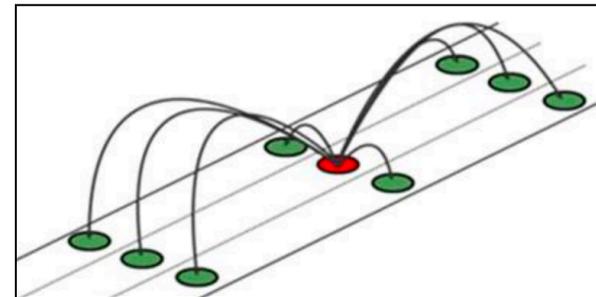
- Given a damaged or incorrectly observed \tilde{x} , the ML system must return an estimate of the original or correct x
 - E.g., ML system is asked to remove dust or scratches from an old photograph



- This requires Multiple outputs
 - Every element of the clean x
 - An understanding of the entire input
 - Since even one damaged area will reveal output as damaged

3. Missing value imputation task

- Given the observations of some elements of \mathbf{x} the model is asked to return estimates of (or probability distribution) over some or all *unobserved* elements of \mathbf{x}

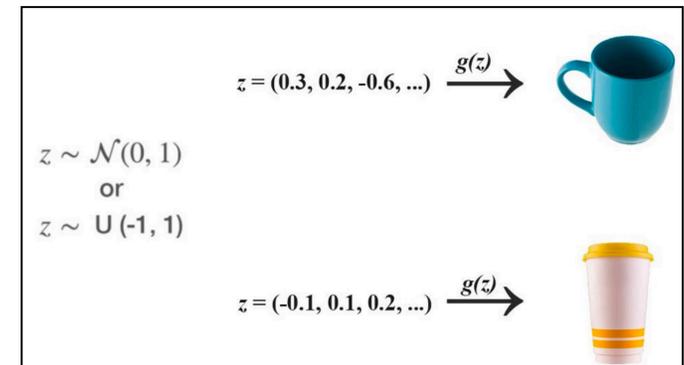


- Or a probability distribution over those parts
- This requires multiple outputs ML algorithm must understand entire input to be able to restore any elements of \mathbf{x}

4. The sampling task

- Model generates samples from $p(\mathbf{x})$
 - Image synthesis

GAN samples noise \mathbf{z} using normal or uniform distribution and utilizes a deep network generator G to create an image \mathbf{x} ($\mathbf{x}=G(\mathbf{z})$)



– Speech synthesis

- Produce waveforms like natural human speech
- Requires multiple output values and model of entire input
- If samples have even one element from wrong distribution, then sampling process is wrong

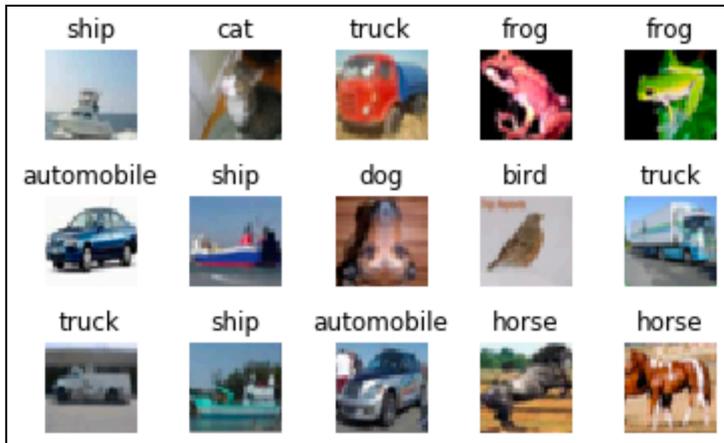
Sampling Task with Natural Images

- Generate new samples from a distribution $p(\mathbf{x})$
 - Probabilistic modeling of natural images

$p(\mathbf{x})$

CIFAR-10 data set

60,000 color images at spatial resolution of 32x32 color (0-255 brightness) from 10 classes



Some Enlarged images from CIFAR-10



Some samples generated by PGM trained on CIFAR-10



Synthesized images displayed here in same order as those that are closest in input set

Model is truly synthesizing new images rather than memorize the training data

Intractability of Rich Distributions

- Computational/statistical challenge
 - Consider CIFAR-10 $32 \times 32 \times 3$ thresholded binary
 - There are 2^{3072} possible images
 - This is 10^{800} larger than no. of atoms in universe
- More generally
 - If \mathbf{x} has n discrete variables with k possible values
 - Representing $P(\mathbf{x})$ by lookup table requires k^n parameters
 - Infeasible for four reasons
 1. Memory
 2. Statistical efficiency
 3. Runtime-inference
 4. Runtime-sampling

Lookup Table

	$P(x=x')$
x^0	0.3
x^1	0.25
x^2	0.27
x^3	0.18

Intractability of Table-based Model

1. *Memory*: cost of storing the representation
 - Unless n and k are small, table would be too large to store
2. *Statistical Efficiency*: As no. of parameters increases, so does training data needed to choose values of parameters using a statistical estimator
 - Astronomically large no of parameters implies astronomically large no of samples to fit accurately

Runtime intractability of naive method

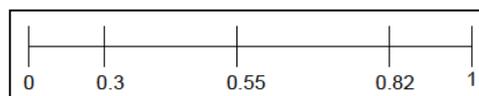
3. Runtime-- Cost of inference

- If we want to use $P(\mathbf{x})$ to compute marginal $P(x_1)$ or conditional $P(x_2|x_1)$
 - It requires summing across entire table
 - Runtime is as high as intractable memory cost of storing model

4. Runtime-- Cost of sampling

- Suppose we need a sample from the model

- Naive way is to draw a sample from $u \sim U(0,1)$



- Iterate through table adding probabilities

- until they exceed u and return outcome corresponding to that position in the table

- Requires reading through entire table--same cost

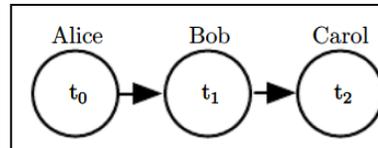
	$P(x=x^i)$
x^0	0.3
x^1	0.25
x^2	0.27
x^3	0.18

Full Table approach is overkill

- Models every possible interaction between variables
- Probability distributions in real tasks are much simpler than this
- Usually variables influence each other only indirectly

Direct and indirect interaction

- Ex: model finishing times in a relay race
- Team has three runners: Alice, Bob, Carol
 - Alice hands baton to Bob, Bob to Carol who finishes



- Alice's finishing time t_0 does not depend on anyone else, Bob's finishing time t_1 depends on Alice's, Carol's t_2 depends on Bob's
 - Carol's finishing time depends only indirectly on Alice's
 - Knowing Bob's finishing time, we are unable to better estimate Carol's finishing time by finding out Alice's finishing time
- Thus model the race using only two interactions:
 1. Alice's effect on Bob
 2. Bob's effect on Carol
 - We can omit the third interaction between Alice and Carol

PGMs model direct interactions

- Allow models to have fewer parameters
 - and therefore be estimated reliably from less data
- Smaller models reduce computational cost of
 - Storing the model
 - Performing inference in the model, and
 - Drawing samples from the model