

Introduction of graphical model, mainly challenge will be how we find an appropriate modeling of variables.

The challenge of Unstructured modeling

- Some applications of probabilistic modeling: Density estimation, denoising, missing value imputation, sampling.
- Modeling all possible structures of input variable (table-based) is not feasible.
- Find a way modeling variables based on their relationships and dependencies, we can get graphical modeling.

Graphical Modeling

- Directed Models (belief network, Bayesian network)

The definition of direction of edges are dependencies. $P(a|b) : \overset{b}{\textcirclearrowleft} \rightarrow \overset{a}{\textcirclearrowright}$

- Local conditional probability contributions

$P(x_i | \text{Par}_i(x_i))$ $P(x_i | X)$ depends on its parents, not on others.

$$\hookrightarrow P(x) = \prod_i P(x_i | \text{Par}_i(x_i))$$

- Conditional Independence means that: we can speculate conditional independence relationships from graphical models.

- Undirected Models (Markov Random Field, Markov Networks)

when we are not clear about the causal relationships, sometimes a mutual influence.

- Clique: The sub-model, aims to model maximum local interactions. for variables within a clique, a potential function $\phi(c)$ is defined, then for the joint probability $\widetilde{P}(x)$ (unnormalized)

$$\widetilde{P}(x) = \prod_{C \in S} \underbrace{\phi(c)}_{\text{Product of all clique potentials}}$$

- The Partition Function

When we get $\tilde{P}(x)$, how we normalize it to be a probability distribution.

$$P(x) = \frac{1}{Z} \tilde{P}(x), \text{ where } Z = \int \tilde{P}(x) dx$$

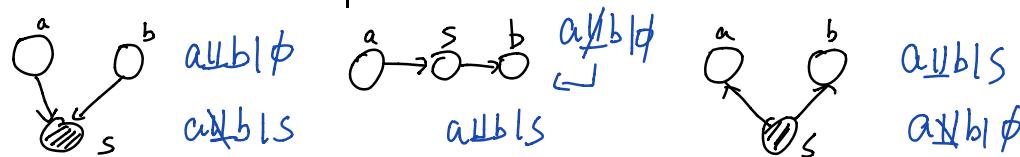
- Sometimes the partition function is intractable, e.g. $Z = \int x^2 dx$, where potential function is a quadratic one. (Domain of x is important)
- For directed graph, modeling is based on probabilistic dependence, thus partition function often do not engage, but for undirected graph, we model interaction of variables using potential function, thus partition function is useful.
- Energy-Based Models

An assumption $\tilde{P}(x) \geq 0, \forall x$ is important, a way of ensuring positive definite of $\tilde{P}(x)$ is energy-based models:

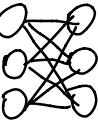
$$\tilde{P}(x) = \exp(-\tilde{\mathcal{E}}(x))$$

Energy function to design
 Ensure positive and convenient of
 calculating log-likelihood.

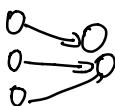
- A term: Boltzmann machine: Considering undirected graph with latent variables.
Markov Random field: " without latent variables.
- Product of cliques can be seen as combining knowledge of experts, each expert may be a clique.
- Separation and D-Separation
 - Undirected graphs: all the paths from a to b is observed (inactive), we can say $a \perp\!\!\!\perp b | s$.
 - Directed graphs: D-separation



- Converting between directed and undirected models

 - Boltzmann machine :  Undirected. latent variable.

Sparse coding :



Directed, coding sparsely

 - Why we need undirected / directed graph respectively ?

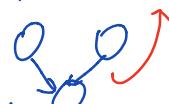
 - How we choose: Which one expresses the most independence with lower representation cost ?

 - How we use: Directed graph \rightarrow Drawing samples

 - Undirected graph \rightarrow Approximate Inference.

 - Converting between them.

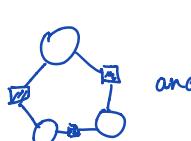
 - Directed \rightarrow Undirected \Rightarrow Dealing with "Immobility" Connections



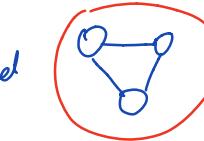
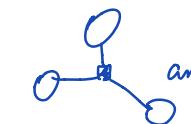
 - Undirected \rightarrow Directed \Rightarrow Adding chord, dealing with loop.

 - Factor Graphs :

 - Solving ambiguity of how the variables actually interact with each other.



and



undirected, no factor

Sampling from Graphical Models

- Ancestral Sampling : Sampling according to the dependence relationship, first of parents and then sampling according to $P(x|P_G(x))$ (Direct Graphical Models)
- Gibbs Sampling : For each time step, sampling only $P(x_i|x_i)$ (Dependence only on neighbours), then iteratively doing sampling till convergence.

Advantage of Structured modeling and learning about dependencies / Inference

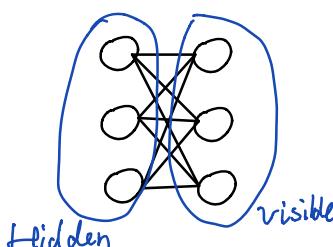
- Advantage: Separation of knowledge representation and inference from knowledge
- Learning dependencies taking use of latent variables: some relationship (dependence) or representations can be modeled efficiently using latent / hidden variables.
- Inference concept: Speculating probabilistic distribution of $P(X | X_{\text{given/observed}})$

Deep learning vs. Structured Model

- Differences:
 - Concept of depth: Computational iterations vs. Hierarchical Dependencies.
 - Representation: Non-explainable, higher dimension of latent variables vs. Explainable, lower dimension of latent variables.
- SGD ←
 - Inference: Distributed representation from $X \rightarrow h$ (approx. inference) vs. Computational feasibility consideration (Sampling, approx. inference, etc.)

Computational feasibility consideration (Sampling, approx. inference, etc.)

- Relation:
 - Tolerance of unknown: Deep learning can be seen as models that have high tolerance of unknown, we focus on representation and gradient.
 - Concepts within them: Graphical models \Rightarrow Latent variables, observed
Deep learning \Rightarrow Hidden Units, input/output.
- An example: RBM seen in two views.



$$P(h|v) = \prod_i P(h_i|v)$$

$$P(v|h) = \prod_i P(v_i|h)$$

Due to restriction property.

Gibbs Sampling: Undirected graph to inference.

Derivatives: $\frac{\partial}{\partial w_{ij}} Z(v, h)$
such as deep learning.

RBM (parametric): learning mapping rules $P(h) = P^*(\text{something})$

Sparse Coding (Non-parametric): learning optimal value with sparse prior.